

Université de Montréal

**Estimateur neuronal de ratio pour l'inférence de la  
constante de Hubble à partir de lentilles  
gravitationnelles fortes**

par

**Ève Campeau-Poirier**

Département de physique  
Faculté des arts et des sciences

Mémoire présenté en vue de l'obtention du grade de  
Maître ès sciences (M.Sc.)  
en physique

21 décembre 2022



# Université de Montréal

Faculté des arts et des sciences

---

Ce mémoire intitulé

## **Estimateur neuronal de ratio pour l'inférence de la constante de Hubble à partir de lentilles gravitationnelles fortes**

présenté par

### **Ève Campeau-Poirier**

a été évalué par un jury composé des personnes suivantes :

*Julie Hlavacek-Larrondo*

---

(présidente-rapporteuse)

*Laurence Perreault Levasseur*

---

(directrice de recherche)

*David Lafrenière*

---

(membre du jury)



## Résumé

---

Les deux méthodes principales pour mesurer la constante de Hubble, soit le taux d'expansion actuel de l'Univers, trouvent des valeurs différentes. L'une d'elle s'appuie lourdement sur le modèle cosmologique aujourd'hui accepté pour décrire le cosmos et l'autre, sur une mesure directe. Le désaccord éveille donc des soupçons sur l'existence d'une nouvelle physique en dehors de ce modèle. Si une autre méthode, indépendante des deux en conflit, soutenait une des deux valeurs, cela orienterait les efforts des cosmologistes pour résoudre la tension.

Les lentilles gravitationnelles fortes comptent parmi les méthodes candidates. Ce phénomène se produit lorsqu'une source lumineuse s'aligne avec un objet massif le long de la ligne de visée d'un télescope. La lumière dévie de sa trajectoire sur plusieurs chemins en traversant l'espace-temps déformé dans le voisinage de la masse, résultant en une image déformée, grossie et amplifiée. Dans le cas d'une source lumineuse ponctuelle, deux ou quatre images se distinguent nettement. Si cette source est aussi variable, une de ses fluctuations apparaît à différents moments sur chaque image, puisque chaque chemin a une longueur différente. Le délai entre les signaux des images dépend intimement de la constante de Hubble.

Or, cette approche fait face à de nombreux défis. D'abord, elle requiert plusieurs jours à des spécialistes pour exécuter la méthode de Monte-Carlo par chaînes de Markov (MCMC) qui évalue les paramètres d'un seul système de lentille à la fois. Avec les détections de milliers de systèmes prévues par l'observatoire Rubin dans les prochaines années, cette approche est inconcevable. Elle introduit aussi des simplifications qui risquent de biaiser l'inférence, ce qui contrevient à l'objectif de jeter la lumière sur le désaccord entre les mesures de la constante de Hubble.

Ce mémoire présente une stratégie basée sur l'inférence par simulations pour remédier à ces problèmes. Plusieurs travaux antérieurs accélèrent la modélisation de la lentille grâce à l'apprentissage automatique. Notre approche complète leurs efforts en entraînant un estimateur neuronal de ratio à déterminer la distribution de la constante de Hubble, et ce, à partir des produits de la modélisation et des mesures de délais. L'estimateur neuronal de ratio s'exécute rapidement et obtient des résultats qui concordent avec ceux de l'analyse traditionnelle sur des simulations simples, qui ont une cohérence statistique acceptable et qui sont non-biaisés.

**Mots-clés:** Constante de Hubble — Lentilles gravitationnelles fortes — Estimateur neuronal de ratio — Cosmologie — Inférence par simulations — Apprentissage automatique.

# Abstract

---

The two main methods to measure the Hubble constant, the current expansion rate of the Universe, find different values. One of them relies heavily on today’s accepted cosmological model describing the cosmos and the other, on a direct measurement. The disagreement thus arouses suspicions about the existence of new physics outside this model. If another method, independent of the two in conflict, supported one of the two values, it would guide cosmologists’ efforts to resolve the tension.

Strong gravitational lensing is among the candidate methods. This phenomenon occurs when a light source aligns with a massive object along a telescope line of sight. When crossing the curved space-time in the vicinity of the mass, the light deviates from its trajectory on several paths, resulting in a distorted and magnified image. In the case of a point light source, two or four images stand out clearly. If this source is also variable, the luminosity fluctuations will appear at different moments on each image because each path has a different length. The time delays between the image signals depend intimately on the Hubble constant.

This approach faces many challenges. First, it requires several days for specialists to perform the Markov Chain Monte-Carlo (MCMC) which evaluates the parameters of a single lensing system at a time. With the detection of thousands of lensing systems forecasted by the Rubin Observatory in the coming years, this method is inconceivable. It also introduces simplifications that risk biasing the inference, which contravenes the objective of shedding light on the discrepancy between the Hubble constant measurements.

This thesis presents a simulation-based inference strategy to address these issues. Several previous studies have accelerated the lens modeling through machine learning. Our approach complements their efforts by training a neural ratio estimator to determine the distribution of the Hubble constant from lens modeling products and time delay measurements. The neural ratio estimator results agree with those of the traditional analysis on simple simulations, have an acceptable statistical consistency, are unbiased, and are obtained significantly faster.

**Keywords:** Hubble constant — Strong gravitational lensing — Neural ratio estimator — Cosmology — Simulation-based inference — Machine learning.





# Table des matières

---

<b>Résumé</b> .....	5
<b>Abstract</b> .....	7
<b>Liste des tableaux</b> .....	13
<b>Liste des figures</b> .....	15
<b>Liste des acronymes</b> .....	17
<b>Liste des symboles</b> .....	21
<b>Remerciements</b> .....	31
<b>Chapitre 1. Introduction</b> .....	33
1.1. Description de ce mémoire .....	34
1.2. Déclaration de l'étudiante .....	34
<b>Chapitre 2. Constante de Hubble</b> .....	37
2.1. Origines .....	37
2.2. Métrique Friedmann-Lemaître-Robertson-Walker .....	39
2.3. Échelle des distances cosmiques .....	40
2.4. Expansion accélérée .....	41
2.5. Crise de la cosmologie .....	42
<b>Chapitre 3. Lentilles gravitationnelles fortes</b> .....	47
3.1. Formalisme .....	47
3.1.1. Angles de déviation .....	48
3.1.2. Équation de lentille .....	49
3.1.3. Potentiel de lentille .....	50
3.2. Applications .....	51

3.2.1.	Observer les sources lumineuses .....	51
3.2.2.	Étudier les défecteurs massifs .....	52
3.2.3.	Inférer la constante de Hubble .....	53
<b>Chapitre 4.</b>	<b>Apprentissage automatique .....</b>	<b>59</b>
4.1.	Perceptron multicouche .....	59
4.2.	Entraînement .....	61
4.2.1.	Fonction de perte .....	62
4.2.2.	Rétropropagation des gradients .....	63
4.2.3.	Descente de gradient .....	64
4.2.4.	Ensembles de données .....	67
4.3.	Inférence par simulations .....	67
4.3.1.	Statistiques bayésiennes .....	68
4.3.2.	Estimateur neuronal de ratio .....	69
4.4.	Set Transformer .....	70
4.4.1.	Attention .....	71
4.4.2.	Attention à têtes multiples .....	71
4.4.3.	Architecture .....	72
<b>Chapitre 5.</b>	<b>Time Delay Cosmography with a Neural Ratio Estimator .....</b>	<b>75</b>
	Résumé .....	76
	Abstract .....	76
5.1.	Introduction .....	77
5.2.	Method .....	78
5.2.1.	Time delay cosmography .....	78
5.2.2.	Neural Ratio Estimator .....	81
5.3.	Simulations .....	82
5.4.	Neural Network .....	84
5.4.1.	Architecture .....	84
5.4.2.	Training .....	85
5.5.	Results .....	86
5.6.	Discussion and Conclusion .....	89

Chapitre 6. Conclusion.....	91
Références bibliographiques .....	95
Annexe A. Neural network variable dimensions .....	101
Annexe B. Congrès où l'étudiante a présenté ses résultats.....	103



## Liste des tableaux

---

5.1	Prior distributions of all the parameters needed to generate Fermat potentials and time delays in our framework .....	83
5.2	Standard deviation of the Gaussian noise distributions used to mimic the uncertainties of lens modeling, time delay measurements, and image position measurements .....	84
A.1	Input sizes for each operation in the NRE.....	101



## Liste des figures

---

2.1	Détermination du coefficient entre la vitesse et la distance des galaxies par Hubble (1929). . . . .	38
2.2	Découverte de l'expansion accélérée de l'Univers par Riess et al. (1998) à l'aide du diagramme de Hubble de supernovae Ia. . . . .	42
2.3	Illustration de la tension entre les déterminations directes et indirectes de $H_0$ par Di Valentino et al. (2021). . . . .	45
3.1	Schéma du formalisme de lentille gravitationnelle . . . . .	50
3.2	Six quasars fortement lentillés observés par le télescope spatial Hubble . . . . .	52
3.3	Détection d'une exoplanète par Beaulieu et al. (2006) grâce à l'effet de microlentille. . . . .	53
3.4	Mesure des délais entre les images d'un quasar fortement lentillé à partir de leurs courbes de lumière par Bonvin et al. (2016). . . . .	55
4.1	Représentation des fonctions d'activation définies par les équations de (4.1.7) à (4.1.3). . . . .	60
4.2	Schéma d'un perceptron multicouche. . . . .	61
4.3	Schéma de la rétropropagation du gradient. . . . .	64
4.4	Représentation de l'effet du momentum par Goodfellow et al. (2016). . . . .	66
4.5	Illustration du mécanisme d'attention et de sa version à têtes multiples par Vaswani et al. (2017) . . . . .	72
5.1	Diagram of gravitational lensing. . . . .	80
5.2	Set Transformer architecture for the discriminator part of the neural ratio estimator . . . . .	85
5.3	Twelve $H_0$ inferences performed by the neural ratio estimator on examples from the test set with zero noise realization. . . . .	87
5.4	Coverage diagnostic of the NRE. . . . .	88

5.5 Population inferences of  $H_0$  with the NRE..... 88



## Liste des acronymes

---

AdaGrad	Gradient adaptatif, de l'anglais <i>Adaptive Gradient algorithm</i>
BNN	Réseau de neurones bayésien, de l'anglais <i>Bayesian Neural Network</i>
CCD	Dispositif à transfert de charges, de l'anglais <i>Charge-Coupled Device</i>
CDM	Matière sombre froide, de l'anglais <i>Cold Dark Matter</i>
CHP	De l'anglais <i>Carnegie Hubble Program</i>
CMB	Fond diffus cosmologique, de l'anglais <i>Cosmic Microwave Background</i>
CRAQ	Centre de Recherche en Astrophysique du Québec
ELU	Unité linéaire exponentielle, de l'anglais <i>Exponential Linear Unit</i>
ESA	Agence spatiale européenne, de l'anglais <i>European Space Agency</i>

FDC	Fond diffus cosmologique
FLRW	Friedmann-Lemaître-Robertson-Walker
GPU	Processeur graphique, de l'anglais <i>Graphics Processing Unit</i>
HPD	Densité a posteriori la plus élevée, de l'anglais <i>Highest Posterior Density</i>
HST	Télescope spatial Hubble, de l'anglais <i>Hubble Space Telescope</i>
iid	Indépendant(e)s et identiquement distribué(e)s
IVADO	Institut de VAlorisation des DONnées
JPL	De l'anglais <i>Jet Propulsion Lab</i>
LMC	Grand Nuage de Magellan, de l'anglais <i>Large Magellanic Cloud</i>
LSST	De l'anglais <i>Legacy Survey of Space and Time</i>
MAB	Bloc d'attention à têtes multiples, de l'anglais <i>Multi-head Attention Block</i>
MCMC	Méthode Monte-Carlo par chaînes de Markov, de l'anglais <i>Markov Chain Monte Carlo</i>

MLCS	Forme de la courbe de lumière multicolore, de l'anglais <i>Multicolore light curve shape</i>
NASA	De l'anglais <i>National Aeronautics and Space Administration</i>
NRE	Estimateur neuronal de ratio, de l'anglais <i>Neural Ratio Estimator</i>
PMA	Échantillonnage par attention à têtes multiples, de l'anglais <i>Pooling by Multi-head Attention</i>
ReLU	Unité linéaire rectifiée, de l'anglais <i>Rectified Linear Unit</i>
rFF	Couche de neurones artificiels, de l'anglais <i>row-wise feedforward layer</i>
RMSProp	Propagation de la racine carrée de l'erreur quadratique moyenne, de l'anglais <i>Root Mean Square Propagation</i>
SAB	Bloc d'auto-attention, de l'anglais <i>Self-Attention Block</i>
SH0ES	De l'anglais <i>Supernovae, <math>H_0</math>, for the Equation of State of dark energy</i>
SIE	Ellipsoïde isotherme singulier, de l'anglais <i>Singular Isothermal Ellipsoid</i>

SN	Supernova
UCLA	Université de Californie à Los Angeles, de l'anglais <i>University of California, Los Angeles</i>
VLA	De l'anglais <i>Very Large Array</i>
WMAP	De l'anglais <i>Wilkinson Microwave Anisotropy Probe</i>

## Liste des symboles

---

$a$	Facteur d'échelle cosmique
$\dot{a}$	Dérivée temporelle du facteur d'échelle cosmique
$a_j$	Activation du neurone artificiel $j$
$A$	Matrice de distorsion
$B$	Nombre d'exemples dans un lot
$\mathcal{B}$	Lot
$c$	Vitesse de la lumière
$c'$	Vitesse de la lumière dans un champ gravitationnel
$\mathcal{C}$	Classe lors d'une tâche de classification
$d$	Dimension des éléments d'un ensemble
$d$	Observation de lentille gravitationnelle

$\mathbf{d}$	Réseau de neurones discriminant
$\mathbf{d}^*$	Réseau de neurones discriminant optimal
$d\ell$	Intervalle infinitésimal du trajet de la lumière
$dr$	Intervalle infinitésimal de rayon spatial
$ds$	Intervalle infinitésimal d'espace-temps
$dt$	Intervalle infinitésimal de temps
$d\vec{x}$	Intervalle infinitésimal d'espace en coordonnées cartésiennes
$d\Omega$	Intervalle infinitésimal spatial d'angle solide
$D$	Dimension d'entrée d'un réseau de neurones artificiels
$D_d$	Distance de diamètre angulaire entre un défecteur massif et un télescope
$D_{ds}$	Distance de diamètre angulaire entre un défecteur massif et une source lumineuse
$D_s$	Distance de diamètre angulaire entre une source lumineuse et un télescope

$D_{\Delta t}$	<i>Time delay distance</i>
$f$	Ratio entre les axes d'un ellipsoïde
$\mathcal{F}(\mathbf{d})$	Intégrant de la fonction de perte d'un NRE
$g$	Cisaillement réduit
$g_{\mu\nu}$	Métrique d'espace-temps
$\mathbf{g}$	Gradient de la fonction de perte par rapport aux poids
$G$	Constante gravitationnelle
$G_{\mu\nu}$	Tenseur d'Einstein
$h$	Ancien symbole du taux d'expansion actuel de l'Univers
h	Nombre de têtes pour l'attention à têtes multiples
$h(\cdot)$	Fonction d'activation
$H$	Nombre de neurones dans une couche cachée
$H(t)$	Paramètre de Hubble

$H_0$	Constante de Hubble
$k$	Courbure de l'Univers
$K$	Coefficient de proportionnalité entre la distance et le décalage vers le rouge
$K$	Dimension de la sortie d'un réseau de neurones artificiels
$\mathcal{K}$	<i>Keys</i>
$l$	Trajet de la lumière
$L$	Résultat d'une normalisation par couche
$\mathcal{L}$	Fonction de perte
$m$	Nombre d'éléments dans un ensemble
$M$	Masse totale
$\mathbf{M}$	Effets observationnels
$n$	Indice de réfraction
$N$	Nombre d'exemples d'entraînement



$O$	Sortie d'une fonction d'attention
$q_0$	Paramètre de décélération de l'Univers
$\mathcal{Q}$	<i>Queries</i>
$r$	Distance de l'origine
$r(\cdot \cdot)$	Ratio de distributions
$\mathbf{r}$	Accumulation de gradient
$R$	Rayon de l'Univers
$\mathbb{R}$	Ensemble des nombres réels
$\mathbf{s}$	Momentum d'ADAMAX
$S$	Vecteur de poids pour l'échantillonnage par attention à têtes multiples
$S_8$	Paramètre cosmologique équivalent à $\sigma_8\sqrt{\Omega_m/0.3}$
$t$	Temps

$t_0$	Âge dynamique de l'Univers
$T_{\mu\nu}$	Tenseur de stress-énergie
$u$	Norme infinie pondérée exponentiellement
$\mathcal{U}$	Distribution uniforme
$\mathbf{v}$	Momentum
$\mathcal{V}$	<i>Values</i>
$\mathbf{w}$	Poids d'un réseau de neurones artificiels
$\vec{x}$	Position spatiale tridimensionnelle
$\mathbf{x}$	Entrée d'un réseau de neurones artificiels
$\mathbf{x}$	Observation
$(x, y)$	Position angulaire dans le champ de vision d'un télescope
$X, Y$	Ensembles entre lesquels l'attention est calculée
$\mathbf{y}$	Cible sur laquelle est entraîné un réseau de neurones artificiels

$\hat{y}$	Sortie d'un réseau de neurones artificiels
$z$	Décalage vers le rouge
$z_d$	Décalage vers le rouge d'une défecteur massif
$z_j$	Unité cachée du neurone artificiel $j$
$z_s$	Décalage vers le rouge d'une source lumineuse
$Z$	Sortie de l'encodeur d'un Set Transformer
$\alpha$	Angle de déviation réduit
$\hat{\alpha}$	Angle de déviation
$\beta$	Position angulaire dans le plan d'une source lumineuse
$\gamma$	Cisaillement
$\delta$	Valeur infinitésimale
$\delta_{ij}$	delta Kronecker
$\Delta$	Différence

$\epsilon$	Paramètre des fonctions d'activation Leaky ReLU et ELU
$\zeta$	Ensemble des paramètres de nuisance
$\eta$	Taux d'apprentissage
$\eta$	Position dans le plan d'une source lumineuse
$\vartheta$	Orientation de la position angulaire dans le plan d'un défecteur massif
$\theta_E$	rayon d'Einstein
$\theta$	Position angulaire dans le plan d'un défecteur massif
$\Theta$	Ensemble des paramètres d'intérêt pour une inférence
$\kappa$	Densité de masse surfacique adimensionnelle
$\lambda$	Constante désignant un membre d'une famille de modèles de masse surfacique dégénérés
$\Lambda$	Constante cosmologique
$\mu$	Amplification

$\xi$	Position dans le plan du défecteur massif
$\pi$	Nombre pi
$\boldsymbol{\pi}$	Ensemble des paramètres d'un modèle
$\rho$	Densité de masse
$\sigma$	Déviatiion standard
$\sigma(\cdot)$	Fonction softmax
$\sigma_8$	Amplitude des fluctuations de densité cosmique
$\Sigma$	Densité de masse surfacique
$\tau$	Époque d'entraînement actuelle
$\phi$	Potentiel de Fermat
$\varphi$	Orientation dans le champ de vision d'un télescope
$\Phi$	Potentiel gravitationnel
$\psi$	Potentiel de lentille

$\omega$	Taux de décroissance
$\Omega_\Lambda$	Densité énergétique du vide
$\Omega_m$	Densité énergétique de la matière
$\Omega_{\text{tot}}$	Densité énergétique totale de l'Univers
$\nabla$	Opérateur de gradient
$\nabla^2$	Opérateur laplacien
$\perp$	Indique les composantes perpendiculaires à un axe
$\odot$	Produit d'Hadamard
$\  \cdot \ $	Norme euclidienne

## Remerciements

---

Merci à ma directrice, Laurence Perreault Levasseur, pour tes réponses, mais aussi pour tes questions, qui m'ont fait progresser. Merci Yashar Hezaveh pour tes conseils sur la recherche et sur la vie.

Merci à Pablo Lemos et particulièrement à Adam Coogan de votre aide à mon projet dans les moments les plus difficiles.

Merci à @ptrblck\_de d'avoir répondu sur des forums à toutes mes questions sur PyTorch, et ce, dix ans avant que je ne me les pose. Merci également à Julia Linhart, la première à avoir porté à mon attention le Set Transformer.

Merci à mes camarades de classe, Alexandre Adam, Jonathan Beaulieu-Émond, Ronan Legin, Laurence Marcotte, Hadi Sotoudeh, Olivier Vincent et Charles Wilson, de votre présence, de votre aide et de votre solidarité, malgré les cours en ligne. Un merci spécial à Charles Wilson de m'avoir accompagnée tout au long de cette aventure et à Alexandre Adam, notre plus dévoué mentor.

Merci à Alexandre Adam, Lucas Goupil, Ana Hoban, Ronan Legin, Olivier Vincent et Charles Wilson pour les soirées de jeux en ligne qui nous ont soudé au cœur du confinement.

Merci à mes formidables compagnons de voyage, Alexandre Adam, Ronan Legin, Hadi Sotoudeh et Charles Wilson, pour les souvenirs inoubliables.

Merci aux membres du comité EDI, Claudia Bielecki, Joseph Choi, Andreas Filipp, Stéphanie Luna, Maria Sadikov, Hadi Sotoudeh, Connor Stone et Benjamin Vigneron, de votre empathie et de votre implication. J'ai hâte de voir où vous mènerez ce projet! Merci aussi aux autres membres du groupe et de l'institut qui ont croisé mon chemin : ça a été un plaisir de travailler avec vous.

J'en profite pour te remercier, Stéphanie Luna, de ta redoutable efficacité en général et de ta détermination indispensable à la réalisation de nos idées.

Merci à Carl Lévesque, Charles Gauthier, Charles Wilson, David Bourbonnais-Sureault, Émile Godbout, Érika Loranger, Eugène Sanscartier, Jonathan Beaulieu-Émond, Laurence

Marcotte, Myriam Prasow-Émond, Pierre-Antoine Bernard, Renaud Girard, Valérie Monette et Vincent Gousy-Leblanc. J'aurais abandonné il y a longtemps si ça n'avait été de vous.

Merci à Alexis Meunier, Camille Huberdeau, Cédric Bisson, Étienne Bernier-Robert, Georges Mainville, Jean-Christophe Cardinal, Justine Faubert, Mathilde Pagé, Mick Frouin et Vincent Williams. Les soirées trop rares en votre compagnie m'ont ressourcée comme rien d'autre.

Merci à Catherine Binette Sirois, Jaëlle L. Landry, Élisabeth Dufresne, Ariane Vallée, Léa-Maude Brodeur et Elisa Zhang pour nos rencontres trop espacées mais toujours rafraîchissantes.

Merci à Sylvain Côté et à Sébastien Wall-Lacelle d'avoir alimenté l'intérêt d'une adolescente de 16-17 ans pour la physique.

Merci à maman, papa, Jérémie et Maxence de votre appui, de vos encouragements et de votre curiosité. Sans même vous en rendre compte, vous avez semé puis entretenu ce qu'il me fallait pour réussir.



# Chapitre 1

---

## Introduction

Six paramètres suffisent à synthétiser toute notre compréhension actuelle de l'Univers à grande échelle. Du moins, ils sont de ce nombre dans le modèle cosmologique largement accepté de nos jours,  $\Lambda$ CDM plat. Par  $\Lambda$ , il inclut l'énergie sombre, d'une origine inconnue, qui s'oppose à la gravitation pour pousser l'Univers à prendre de l'expansion de plus en plus rapidement. La matière sombre froide (CDM), aussi d'une nature indéterminée, y représente la vaste majorité de l'énergie sous forme de matière et dicte la configuration des grandes structures de l'Univers. Sa platitude ne tient qu'à son espace-temps sans courbure en l'absence de masse pour le déformer.

Le modèle  $\Lambda$ CDM réussit brillamment à rendre compte des grandes échelles de l'Univers et de son évolution dans le temps. En revanche, certaines failles laissent croire qu'il est incomplet. Par exemple, il surestime la densité et les cuspidés des noyaux galactiques de matière sombre (Flores & Primack, 1994; Moore, 1994). Les observations indiquent aussi que le nombre de petites galaxies et de satellites nains, ainsi que leurs masses, sont inférieurs aux prédictions (Klypin et al., 1999; Moore et al., 1999). Par ailleurs, une tension persiste entre différentes mesures de deux de ses paramètres, soit l'amplitude des structures cosmiques ( $S_8$ ) et la constante de Hubble ( $H_0$ ; Collaboration Planck et al., 2014).

Pour cette dernière, la mesure de la collaboration Planck, effectuée à partir du spectre de puissance du fond diffus cosmologique (FDC), contrevient à celle obtenue de façon directe par la collaboration *Supernovae,  $H_0$ , for the Equation of State of dark energy* (SH0ES), c'est-à-dire par la mise en relation du décalage vers le rouge et la distance de corps célestes. Or, la méthode de la collaboration Planck est intimement liée au modèle  $\Lambda$ CDM. Si sa valeur s'avère inexacte, cela suggérerait l'existence d'un phénomène physique dérogeant à  $\Lambda$ CDM. Toutefois, une erreur expérimentale, d'un côté comme de l'autre, pourrait tout aussi bien engendrer ce désaccord. Pour trancher entre ces deux explications, une autre méthode, indépendante du FDC et des distances cosmiques ainsi que suffisamment précise et robuste, doit prendre parti.

Le recours aux lentilles gravitationnelles fortes répond au premier critère. Le *Legacy Survey of Space and Time* (LSST) s'inscrit dans les efforts déployés pour remplir le second. Au cours des 10 prochaines années, ce relevé astronomique anticipe des milliers d'observations de lentilles gravitationnelles propices à l'inférence de  $H_0$ , ce qui devrait améliorer considérablement la précision et l'exactitude des résultats (Oguri & Marshall, 2010). En revanche, l'analyse d'un seul de ces systèmes requiert plusieurs jours et demande régulièrement l'attention des spécialistes. La taille des données recueillies par le LSST impose une accélération de cette analyse. De plus, la méthode actuelle introduit des assertions simplistes qui plombent sa crédibilité.

Dans ce mémoire, je développe un procédé rapide et flexible pour mesurer la constante de Hubble à partir des lentilles gravitationnelles fortes. Pour ce faire, j'entraîne et je teste sur des simulations un estimateur neuronal qui évalue la distribution de  $H_0$  en fonction d'observables et de produits d'analyse. Mon travail remplace donc la méthode longue et fastidieuse de Monte-Carlo par chaînes de Markov (MCMC) habituellement appliquée aux paramètres du système de lentille. Les résultats de l'estimateur s'obtiennent rapidement, ils concordent avec ceux de l'analyse traditionnelle et sont non-biaisés, mais présentent des incertitudes faiblement surestimées. Ce mémoire amène ainsi la stratégie des lentilles gravitationnelles plus près de la vitesse, de la précision et de la fiabilité nécessaires à la résolution de la crise en cosmologie.

## 1.1. Description de ce mémoire

Ce mémoire vise à démontrer la faisabilité, la rapidité et l'exactitude d'une inférence par simulations, basée sur l'apprentissage automatique, pour déterminer la constante de Hubble à partir de lentilles gravitationnelles fortes.

Le chapitre 2 résume l'historique des mesures de  $H_0$  et dresse le portrait de la crise actuelle. Le chapitre 3 détaille le formalisme et les applications des lentilles gravitationnelles fortes, dont la détermination de  $H_0$ . Le chapitre 4 introduit d'abord certaines bases de l'apprentissage automatique, soient les réseaux de neurones artificiels et leur entraînement. Ses sections 4.3 et 4.4 se consacrent respectivement à l'inférence par simulations et au Set Transformer qui interviennent dans le chapitre 5. Ce dernier rapporte les résultats d'une inférence de  $H_0$ , réalisée par un réseau de neurones artificiels, sur des simulations de lentilles gravitationnelles fortes.

## 1.2. Déclaration de l'étudiante

Je, Ève Campeau-Poirier, déclare que l'entièreté du travail présenté dans ce mémoire est le mien. Les articles et autres ouvrages de référence qui y sont cités proviennent de ma

propre revue de littérature. Les sources des figures que je n'ai pas conçues moi-même sont inscrites dans leur description.

J'ai rédigé l'article au chapitre 5, mais le crédit des idées scientifiques revient à Laurence Perreault Levasseur et à Yashar Hezaveh. J'ai écrit l'entièreté du code des simulations, excepté un emprunt aux notes de cours de Massimo Meneghetti pour résoudre l'équation de lentille. J'ai implémenté le réseau de neurones à partir de l'exemple fournit par Juho Lee <sup>1</sup>. J'ai effectué ses entraînements et les recherches d'hyperparamètres. Finalement, j'ai produit et analysé les résultats en bénéficiant des conseils techniques d'Adam Coogan.

---

<sup>1</sup>[https://github.com/juho-lee/set\\_transformer](https://github.com/juho-lee/set_transformer)



# Chapitre 2

---

## Constante de Hubble

*Nothing compares to the measurement of the Hubble constant in bringing out the worst in astronomers.*

---

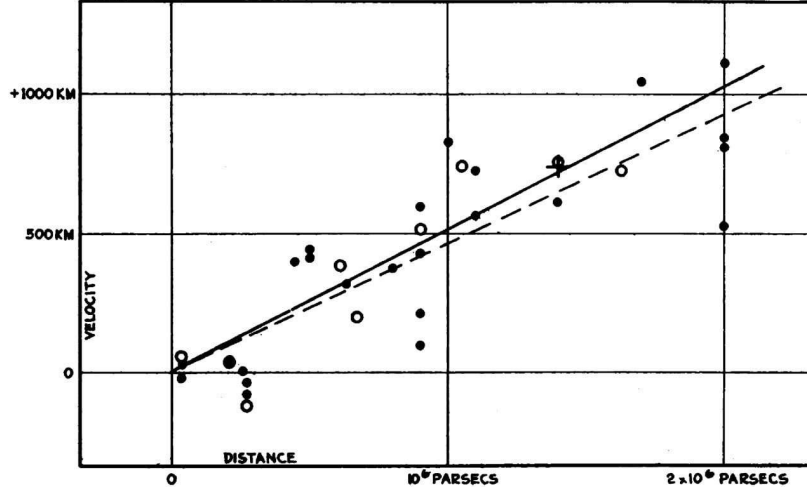
Christopher S. Kochanek (2006)

### 2.1. Origines

Vesto Slipher (1917) constate que la majorité des galaxies, appelées nébuleuses extra-galactiques à l'époque, s'éloignent de la Terre. Edwin Hubble (1929) se propose de déterminer le coefficient  $K$ , soit le facteur de proportionnalité entre la vitesse et la distance des nébuleuses extra-galactiques par rapport à la Terre, et ce, à partir des données les plus fiables, puisque les tentatives précédentes s'étaient avérées peu concluantes. Les mesures de distance provenaient de la luminosité apparente d'étoiles dont la luminosité absolue était connue, c'est-à-dire des « chandelles standards » dont font partie les céphéides. Les vitesses avaient été évaluées grâce au décalage des raies spectrales  $z$  par effet Doppler. Hubble estime  $K$  autour de  $500 \text{ km s}^{-1} \text{ Mpc}^{-1}$  (voir Fig. 2.1) et attribue la corrélation à l'effet de Sitter, soit au ralentissement des vibrations atomiques et à la tendance des particules de matière à se disperser.

Deux ans plus tôt, Georges Lemaître (1927) avait conclu que le rayon de l'Univers croît. Il avait aussi suggéré que cette expansion causait l'éloignement des nébuleuses extra-galactiques observé par Slipher. Du même coup, Lemaître expliquait aussi la corrélation entre la vitesse de ces corps célestes et leur distance, corrélation précédemment étudiée par Hermann Weyl (1923), Arthur Eddington (1923), Ludwik Silberstein (1924) et Knut Lundmark (1924, 1925), notamment. Surtout, utilisant les vitesses et les distances de 42 nébuleuses mesurées par Hubble et Strömberg, en plus d'assumer une relation linéaire entre ces variables, Lemaître avait déterminé que la vitesse moyenne de ces corps célestes à 1 Mpc était de  $625 \text{ km s}^{-1}$ .

Cette mesure, équivalente au coefficient  $K$ , devance celle de Hubble. De plus, l'explication du



**Fig. 2.1.** Détermination du coefficient entre la vitesse et la distance des galaxies. Les distances sont obtenues par la luminosité apparente des céphéides, tandis que les vitesses sont déterminées par effet Doppler. La ligne pleine et les cercles noirs considèrent les galaxies individuellement. La ligne pointillée et les cercles blancs regroupent les galaxies. La croix indique la vitesse moyenne correspondant à la distance moyenne de 22 galaxies dont les distances individuelles ne pouvaient être estimées. Crédit : Hubble (1929)

phénomène proposée par Lemaître est celle qui s'avèrera exacte. Néanmoins, tandis que les travaux de Lemaître passent initialement sous le radar, ceux de Hubble deviennent la référence principale concernant la corrélation entre la distance des nébuleuses extra-galactiques et leur vitesse. Les termes pour nommer cette corrélation dans la littérature sont « loi de Hubble » (Milne, 1933), « facteur de Hubble » (Haas, 1938), puis « constante de Hubble » (Behr, 1951) qui désigne encore aujourd'hui le taux d'expansion actuel de l'Univers. Aussi, Robertson (1933) introduit le symbole  $h$ , en référence à Hubble, pour représenter le taux d'expansion de l'Univers. Il l'écrit comme :

$$h \equiv (R'/R)_0 , \quad (2.1.1)$$

où  $R$  est le rayon de l'Univers, où le prime représente une dérivée temporelle et où l'indice 0 indique l'époque actuelle, ce qui est très près de la définition présente de  $H_0$ ,

$$H_0 \equiv \frac{\dot{a}(t_0)}{a(t_0)} , \quad (2.1.2)$$

où  $t$  est le temps,  $a$ , le facteur d'échelle et  $\dot{a}$ , la dérivée temporelle de ce dernier.

Par ailleurs, des alternatives à l'expansion de l'Univers pour expliquer la corrélation entre la distance et le décalage vers le rouge perdureront longtemps. Par exemple, Fritz Zwicky (1929) propose une perte d'énergie des photons due à leur voyage, phénomène nommé « lumière fatiguée ». Le test suivant permet de discriminer cette hypothèse de celle de Lemaître : si l'Univers est en expansion, un facteur  $(1 + z)$  de dilatation temporelle affecte le signal

d'événements lointains, tandis que si la lumière se fatigue, le signal est simplement atténué. Il faut attendre que Perlmutter et al. (1995) montrent que les courbes de lumière des supernovae subissent effectivement une dilatation temporelle pour éliminer l'hypothèse de la lumière fatiguée.

## 2.2. Métrique Friedmann-Lemaître-Robertson-Walker

Malgré tout, le concept d'Univers dynamique s'implante en cosmologie et la constante de Hubble s'intègre éventuellement à son cadre théorique. Avant Lemaître, Alexander Friedmann (1922, 1924) avait publié la première solution à l'équation de champ d'Einstein qui n'assumait pas un univers statique. Arthur Geoffrey Walker (1937) et Howard P. Robertson (1935, 1936a,b) poursuivent les travaux en décrivant un espace-temps spatialement homogène et isotrope, mais variable en fonction du temps. Le fruit de toutes ces contributions est la métrique Friedmann-Lemaître-Robertson-Walker (FLRW),

$$ds^2 = -dt^2 + \frac{a^2(t)dr^2}{a^2(t) - kr^2} + r^2d\Omega^2, \quad (2.2.1)$$

où  $ds$  est un élément d'espace-temps,  $dt$  est un élément de temps,  $dr$  est un élément spatial de rayon,  $d\Omega$  est un élément spatial d'angle solide,  $a(t)$  est une fois de plus le facteur d'échelle qui représente le dynamisme de l'Univers et  $k$  est la courbure de l'espace-temps.

La théorie de la relativité générale développée par Albert Einstein (1915) stipule que tout contenu énergétique affecte la géométrie de l'espace-temps. L'équation qui régit ce principe est l'équation de champ d'Einstein,

$$G_{\mu\nu} + \Lambda g_{\mu\nu} = 8\pi T_{\mu\nu}, \quad (2.2.2)$$

où  $G_{\mu\nu}$  est le tenseur d'Einstein défini par la métrique,  $\Lambda$  est la constante cosmologique associée à l'énergie du vide,  $g_{\mu\nu}$  est la métrique et  $T_{\mu\nu}$  est le tenseur de stress-énergie.

À partir de cette équation et de la métrique FLRW, on trouve que la courbure de l'Univers  $k$  obéit à :

$$-k = a^2(t)H^2(t)(1 - \Omega_{\text{tot}}), \quad (2.2.3)$$

où  $H$  est le paramètre de Hubble,

$$H(t) \equiv \frac{\dot{a}(t)}{a(t)}. \quad (2.2.4)$$

L'évaluation du paramètre de Hubble à l'époque actuelle équivaut à  $H_0$ . Quant à  $\Omega_{\text{tot}}$ , il indique la densité d'énergie totale de l'Univers, normalisée par rapport à la valeur critique pour laquelle  $k = 0$ .

Bien que  $H_0$  s'imbrique adéquatément dans la description mathématique de l'Univers ainsi étoffée durant les années 30, sa mesure demeure sujet à controverse.

## 2.3. Échelle des distances cosmiques

Puisque leurs données provenaient de galaxies locales, Hubble et de Lemaître ont surestimé  $H_0$ . À ces distances, les vitesses particulières dominent celle due à l'expansion de l'Univers, ce qui empêche d'évaluer cette dernière adéquatement. Or, la correction de ce biais ne scelle pas la valeur de la constante de Hubble. De la fin des années 50 à la fin des années 90, une dissension persiste entre les valeurs faibles de  $H_0$ , principalement mesurées par Allan Sandage et Gustav Andreas Tammann, et les valeurs élevées, défendues entre autres par Gérard de Vaucouleurs et Sidney van den Bergh. Les valeurs faibles se maintiennent autour de  $50 \text{ km s}^{-1} \text{ Mpc}^{-1}$  à partir de 1975. Celles élevées couvrent initialement un intervalle entre  $100$  et  $200 \text{ km s}^{-1} \text{ Mpc}^{-1}$ , mais elles diminuent graduellement pour se situer entre  $70$  et  $90 \text{ km s}^{-1} \text{ Mpc}^{-1}$  (Trimble, 1996). Au début des années 90, le noyau du problème réside dans l'échelle des distances cosmiques (Fukugita et al., 1993).

Cette dernière provient du fait que plusieurs méthodes permettent de mesurer les distances cosmiques, mais qu'elles ne fonctionnent toutes qu'à une échelle spécifique. Pour évaluer la distance d'un objet lointain, il faut donc déterminer la distance relative entre de multiples « échelons » jusqu'à cet objet, et ce, par différentes techniques. Par exemple, la luminosité absolue des céphéides et des étoiles Lyrae RR dépend intimement de leur période lumineuse. Cette dernière, combinée à leur luminosité apparente, fournit donc une estimation de leur distance et, conséquemment, de celle de leur galaxie hôte. Or, elles ne brillent pas assez fort pour être détectées dans des galaxies lointaines. Un autre moyen est alors requis pour déterminer la distance entre les galaxies locales et celles plus éloignées.

Les valeurs de  $H_0$  autour de  $50 \text{ km s}^{-1} \text{ Mpc}^{-1}$  résultent de l'emploi des supernovae de type Ia (SNe Ia) comme chandelle standard. Dans les années 60, des modèles théoriques identifient les supernovae comme des explosions thermonucléaires de naines blanches, majoritairement composées de carbone et d'oxygène, ayant atteint la masse de Chandrasekhar (Hoyle & Fowler, 1960; Arnett, 1969; Colgate & McKee, 1972). Kowal (1968) montre que les SNe de type I suivent un diagramme de Hubble bien défini, et donc, qu'elles pourraient servir à l'inférence de  $H_0$ . Aussi, Phillips (1993) découvre que leur chute de brillance dans les 15 jours suivant le pic trahit leur luminosité intrinsèque. La distance d'une galaxie locale, hôte d'une SNe Ia, se mesure grâce à ses céphéides. La distance et la luminosité de cette SNe, maintenant connues, calibrent celles des SNe Ia lointaines, d'où leur utilisation comme chandelle standard.

Les déterminations de  $H_0$  autour de  $80 \text{ km s}^{-1} \text{ Mpc}^{-1}$  avaient plutôt recours à la relation de Tully-Fisher (Tully & Fisher, 1977), aux fluctuations de brillance superficielle (Tonry, 1991) et à la fonction de luminosité des nébuleuses planétaires (Jacoby et al., 1990) pour dépasser les distances mesurables par l'observation des céphéides. La relation de Tully-Fisher lie la luminosité d'une galaxie spirale à sa vitesse de rotation. Cette dernière s'évalue par l'effet



Doppler qui décale le spectre de la galaxie. La deuxième méthode consiste à estimer la distribution de luminosité des étoiles d’une galaxie à partir de la variance statistique de la brillance superficielles de ses comparses. Quant à la fonction de luminosité des nébuleuses planétaires, elle exhibe une nette limite supérieure. Des preuves empiriques indiquent que cette limite est universelle, ce qui en fait un moyen de calibration des luminosités et des distances. Ces trois méthodes présentent une bonne précision et s’accordent les unes avec les autres, mais entrent en conflit avec les SNe Ia.

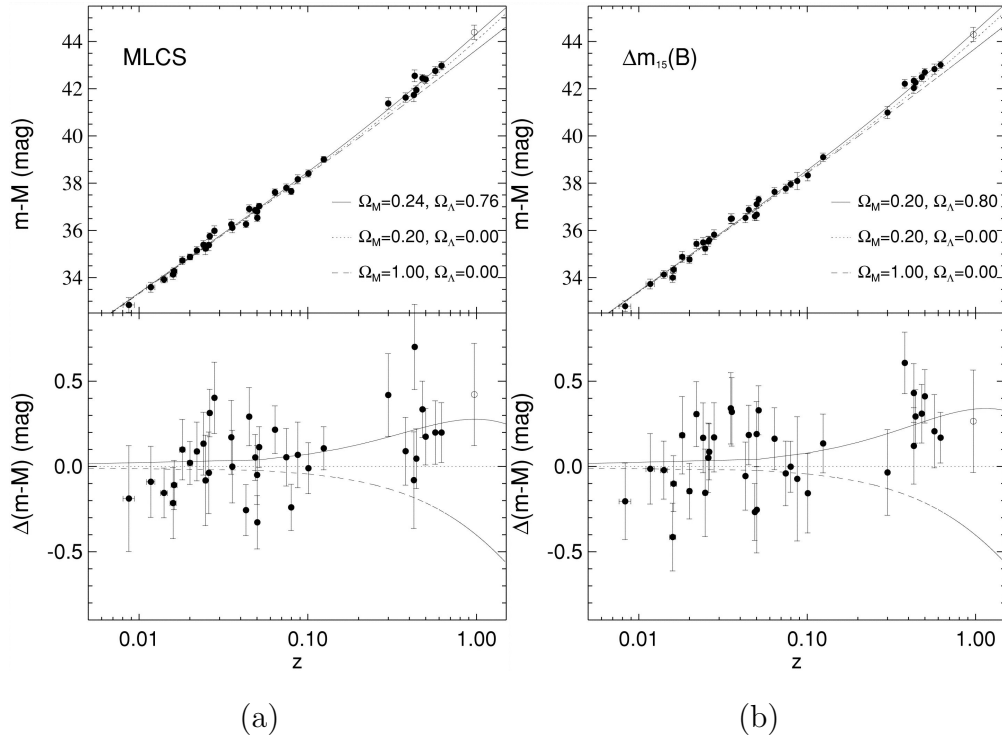
## 2.4. Expansion accélérée

L’étude des SNe Ia par le relevé astronomique Calán/Tololo (Hamuy et al., 1993) améliore significativement la compréhension de celles fortement décalées vers le rouge, et donc, remédie à certaines erreurs systématiques. Ces dernières sont affectent l’ensemble des mesures de façon similaire, tandis que les erreurs statistiques varient d’une mesure à l’autre. Les données de Calán/Tololo ont notamment mené au développement d’une méthode de mesure des distances basée sur la courbe de lumière des SNe Ia en quatre couleurs photométriques (MLCS) (Riess et al., 1996). Le *Supernova Cosmology Project* (Perlmutter et al., 1995) et *High- $z$  Supernova Search Team* (Schmidt et al., 1998) se lancent alors dans la recherche et l’analyse de lointaines SNe Ia.

Adam Riess et al. (1998) apportent ainsi un nouvel élément au débat. Leur étude se penche sur des observations spectroscopiques et photométriques de 50 supernovae de type Ia, dont 16 fortement décalées vers le rouge et 34 locales. L’objectif consiste à imposer des limites aux valeurs de  $H_0$ , de la densité énergétique de la matière  $\Omega_m$ , de la densité énergétique du vide  $\Omega_\Lambda$ , du paramètre de décélération  $q_0$  et de l’âge dynamique de l’Univers  $t_0$ . Différentes méthodes d’ajustement des paramètres, différents sous-ensembles de données et différentes distributions a priori favorisent unanimement  $\Omega_\Lambda > 0$  et  $q_0 < 0$ , c’est-à-dire une accélération de l’expansion de l’Univers (voir Fig. 2.2). Perlmutter et al. (1999) corroborent ce résultat.

Jusqu’alors,  $\Lambda = 0$  était largement accepté par la communauté scientifique. Ces publications marquent donc l’avènement d’un nouveau modèle cosmologique,  $\Lambda$ CDM. Dans ce modèle, l’énergie de l’Univers est dominée par  $\Lambda$  qui représente maintenant une énergie dite « sombre », accélérant l’expansion de l’Univers. L’autre contenu énergétique d’importance est la matière, dont la forme la plus abondante est dite « sombre et froide », c’est-à-dire de nature inconnue et n’interagissant que gravitationnellement. Selon les dernières mesures de Collaboration Planck et al. (2020), leurs densités sont respectivement  $\Omega_\Lambda = 0,6847 \pm 0,0073$  et  $\Omega_m = 0,3153 \pm 0,0073$ , ce qui se traduit par un Univers plat.

Suite à cette découverte, plusieurs équipes, se fiant à différentes échelles de distance, obtiennent des valeurs en accord. Jha et al. (1999) évaluent  $H_0$  à  $68_{-6}^{+8}$  km s<sup>-1</sup> Mpc<sup>-1</sup> et Mould



**Fig. 2.2.** Le graphique du haut montre le diagramme de Hubble de SNe Ia. Les distances sont mesurées par (a) la méthode MLCS et (b) par ajustement d'un modèle à la courbe de lumière. Pour un univers plat, les paramètres qui s'accordent le mieux avec les données sont (a)  $\Omega_m = 0,24$  et  $\Omega_\Lambda = 0,76$  ainsi que (b)  $\Omega_m = 0,20$  et  $\Omega_\Lambda = 0,80$ . Le graphique du bas compare les données au modèle  $\Omega_m = 0,20$  et  $\Omega_\Lambda = 0$ . Le cercle blanc indique une SN sans classification spectroscopique et sans mesure de couleur. Crédit : Riess et al. (1998)

et al. (2000), à  $71 \pm 6 \text{ km s}^{-1} \text{ Mpc}^{-1}$ . Pour le *Key Project* du télescope spatial Hubble (HST), Freedman et al. (2001) trouvent  $72 \pm 8 \text{ km s}^{-1} \text{ Mpc}^{-1}$ . Sandage et al. (2006) mesurent  $62,3 \pm 1,3$  (aléatoire)  $\pm 5,0$  (systématique)  $\text{km s}^{-1} \text{ Mpc}^{-1}$ . De plus, la collaboration de la *Wilkinson Microwave Anisotropy Probe* (WMAP), qui recueille les spectres de puissance angulaires de la température et de la polarisation du FDC, publie un résultat cohérent avec ceux basés sur le diagramme de Hubble. Spergel et al. (2003) rapportent  $H_0 = 72 \pm 5 \text{ km s}^{-1} \text{ Mpc}^{-1}$  en ajustant les paramètres d'un modèle d'univers plat et dominé par  $\Lambda$  aux données du FDC. Les incertitudes laissent toutefois place à l'amélioration.

## 2.5. Crise de la cosmologie

La publication des résultats de Collaboration Planck et al. (2014) déclenche un nouveau désaccord. En ajustant les six paramètres du modèle cosmologique  $\Lambda$ CDM plat aux données du FDC, la collaboration détermine que  $H_0 = 67,3 \pm 1,2 \text{ km s}^{-1} \text{ Mpc}^{-1}$ . Or, les dernières mesures basées sur les céphéides et les supernovae, effectuées par Riess et al. (2011) pour le

programme *Supernovae,  $H_0$ , for the Equation of State of Dark energy* (SH0ES) et par Freedman et al. (2012) pour le *Carnegie Hubble Program* (CHP), sont respectivement  $73,8 \pm 2,4$  km s<sup>-1</sup> Mpc<sup>-1</sup> et  $74,3 \pm 2,1$  (systématique) km s<sup>-1</sup> Mpc<sup>-1</sup>. La tension entre les estimations de  $H_0$  provenant du FDC et des SNe Ia s'élève alors à  $2,5 \sigma$ .

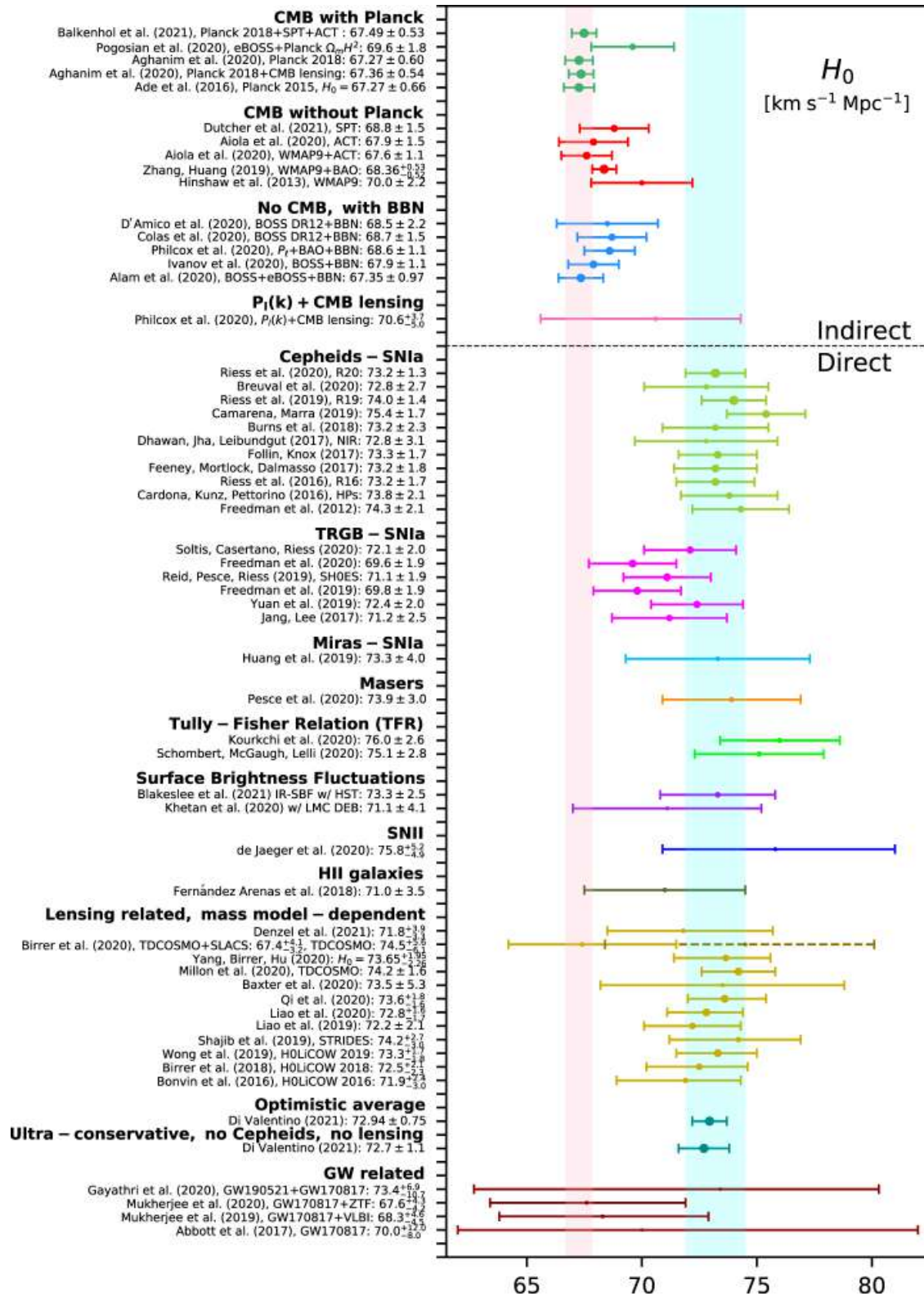
Comme la méthode de la collaboration Planck dépend fortement du modèle cosmologique, cette divergence indique que  $\Lambda$ CDM est possiblement incomplet. Collaboration Planck et al. (2014) notent, par exemple, qu'une extension à la physique des neutrinos ou que la dynamique de l'énergie sombre mitigerait l'écart entre les mesures, mais que les données du FDC n'appuient fermement ni l'une ni l'autre. Par ailleurs, la collaboration remarque des disparités entre les prédictions de  $\Lambda$ CDM et les ordres multipolaires les plus bas du spectre de puissance du FDC, même s'ils concordent rigoureusement aux ordres plus élevés. Aussi, elle fait valoir que le manque de connaissances sur les premières phases de l'Univers laissent place à de la nouvelle physique. Des erreurs systématiques qui n'ont pas été prises en compte engendreraient également un tel désaccord. Pour l'inférence de Riess et al. (2011), l'incertitude provient majoritairement du premier échelon de l'échelle des distances cosmiques. Collaboration Planck et al. (2014) soulignent qu'une légère sous-estimation de cette incertitude suffirait à expliquer l'incompatibilité des valeurs obtenues à partir du FDC et des supernovae.

La méthode de SH0ES repose sur les observations prises par le HST des céphéides qui se trouvent dans les galaxies hôtes des SNe Ia pour calibrer la relation entre le décalage vers le rouge et la magnitude apparente de celles-ci. Ses échelons de distances comprennent la Voix Lactée, la galaxie d'Andromède (M31), la galaxie spirale NGC 4258 et la galaxie satellite du Grand Nuage de Magellan (LMC). Riess et al. (2016, 2018, 2019, 2022) entreprennent une vérification exhaustive des incertitudes et des erreurs systématiques liées à cette approche. Ils observent de nouvelles céphéides dans les galaxies hôtes des SNe Ia (en proche infrarouge), dans le LMC et dans M31. Ils mesurent à nouveau les parallaxes des céphéides de la Voix Lactée avec HST et Gaia. Ils exploitent différents phénomènes astrophysiques pour améliorer la précision sur les distances de NGC 4258 et du LMC. Ces articles étudient l'effet de la période de luminosité, de la poussière et de la métallicité des céphéides. Ils analysent aussi l'impact sur l'inférence du choix des échelons, des relevés astronomiques utilisés, de l'étendue des décalages vers le rouge, de la couleur des SNe Ia et de la correction des vitesses particulières.

Durant la même période, la collaboration Planck se dédie également à la recherche d'erreurs systématiques. Collaboration Planck et al. (2016) constatent qu'un artefact dans le spectre de puissance du FDC, causé par de l'interférence entre deux détecteurs, n'avait pas été adéquatement retiré. L'impact sur les paramètres cosmologiques inférés est toutefois minime.

Collaboration Planck et al. (2020) décèlent aussi des erreurs dans les efficacités de polarisation, dont la correction entraîne de meilleures contraintes. En effet, tous ces efforts des deux côtés ne mènent qu'à une réduction de leurs incertitudes respectives, aggravant la tension à  $5\sigma$  entre les valeurs les plus récentes de Planck et de SH0ES, respectivement de  $67,4 \pm 0,5$  km s<sup>-1</sup> Mpc<sup>-1</sup> (Collaboration Planck et al., 2020) et de  $73,04 \pm 1,04$  km s<sup>-1</sup> Mpc<sup>-1</sup> (Riess et al., 2022).

De plus, la dichotomie ne se limite plus à ces deux méthodes : elle englobe maintenant deux catégories. La première rassemble les mesures directes, provenant d'observations de l'Univers local et indépendantes du modèle cosmologique, telles que le diagramme de Hubble, les lentilles gravitationnelles et les ondes gravitationnelles. La seconde correspond aux mesures indirectes, issues de l'Univers lointain et dépendantes du modèle cosmologique, dont le fond diffus cosmologique, avec ou sans les données de Planck, et la nucléosynthèse primordiale (voir Fig. 2.3). Ces résultats expérimentaux semblent pointer vers l'option d'un modèle cosmologique incomplet. Pourtant, selon Schöneberg et al. (2022), les modifications au modèle  $\Lambda$ CDM proposées pour résoudre la tension sur  $H_0$  n'améliorent pas, voire empirent, celle sur  $S_8$ , un autre paramètre cosmologique qui quantifie l'agglomération de la matière. En effet, un nouveau modèle cosmologique devrait prédire des valeurs plus élevées que  $\Lambda$ CDM pour ces deux paramètres afin de remédier aux tensions avec leurs mesures locales. Or, augmenter le taux d'expansion de l'Univers a un effet dispersif qui diminue l'agglomération de la matière. Le mystère reste donc entier.



**Fig. 2.3.** Comparaison entre les déterminations directes et indirectes de  $H_0$ . La bande rose correspond à la mesure de Collaboration Planck et al. (2020) et la cyan, à celle de Riess et al. (2021). Crédit : Di Valentino et al. (2021)



# Chapitre 3

---

## Lentilles gravitationnelles fortes

Une lentille gravitationnelle se produit lorsqu'une source lumineuse s'aligne en arrière-plan d'un objet massif le long d'une ligne de visée. Les sources de lumière d'intérêt sont les galaxies et les quasars. Quant à l'objet massif, il s'agit le plus souvent d'une galaxie ou d'un amas de galaxies. Les rayons lumineux traversant l'espace-temps courbé par l'objet massif dévient de leur trajectoire initiale. L'image détectée de la source lumineuse s'en trouve déformée, mais aussi magnifiée, ce qui rappelle l'action d'une lentille.

Si l'alignement est insuffisant, l'image sera seulement étirée ou cisailée, ce qu'on qualifie de lentille gravitationnelle « faible ». Si l'alignement est adéquat, l'image est considérablement distordue. La lentille gravitationnelle, dite « forte » dans ce cas, génère plusieurs images en redirigeant les photons sur plus d'un trajet. Si la source lumineuse est ponctuelle, typiquement deux ou quatre images se distinguent. Si la source est étendue, l'image prend plutôt la forme d'arcs ou d'un anneau lumineux.

Les millilentilles et les microlentilles sont des déclinaisons à petites échelles des lentilles gravitationnelles fortes. Dans ces cas, un objet de masse stellaire fait office de défecteur. La proximité des images qu'il produit ne permet pas aux télescopes de discerner chacune d'elles, mais la variation de luminosité est perceptible. Les milli et microlentilles se remarquent surtout lorsqu'elles jouent le rôle de seconde lentille, c'est-à-dire lorsqu'elles amplifient une image déjà magnifiée par une lentille gravitationnelle de grande taille.

### 3.1. Formalisme

Le formalisme suivant provient de l'ouvrage de référence de Bartelmann & Schneider (2001) et des notes de cours de Meneghetti. Il assume que les dimensions du défecteur sont inférieures à celle du système optique, c'est-à-dire aux distances qui séparent la source, le défecteur et l'objectif. De plus, il considère que le potentiel gravitationnel Newtonien  $\Phi$  du défecteur est faible, c'est-à-dire qu'il satisfait  $\Phi/c^2 \ll 1$ . Cette condition tient dans

le voisinage des galaxies et signifie que les rayons dévieront de leur trajectoire à cause du potentiel, mais qu'ils pourront y échapper.

### 3.1.1. Angles de déviation

Un tel déflecteur perturbe la métrique de l'espace-temps de Minkowski, engendrant un élément de ligne,

$$ds^2 = \left(1 + \frac{2\Phi}{c^2}\right) c^2 dt^2 - \left(1 - \frac{2\Phi}{c^2}\right) (d\vec{x})^2, \quad (3.1.1)$$

où  $c$  est la vitesse de la lumière,  $t$  est la dimension temporelle et  $\vec{x}$  représente les dimensions spatiales.

Le principe de Fermat stipule que la lumière emprunte des trajets dont le temps de voyage,

$$t = \int \frac{n(\vec{x}(l))}{c} dl, \quad (3.1.2)$$

est un extremum. L'indice de réfraction  $n$  dépend de la position dans l'espace  $\vec{x}$ , laquelle est paramétrisée par  $l$  le long du trajet de la lumière. Il correspond ici au ratio entre la vitesse de la lumière dans le vide et celle dans le champ gravitationnel du déflecteur. À partir de la métrique (3.1.1) et de la supposition que  $\Phi/c^2 \ll 1$ , la vitesse de la lumière dans le champ gravitationnel du déflecteur est

$$c' = \frac{\|d\vec{x}\|}{dt} = c \sqrt{\frac{1 + \frac{2\Phi}{c^2}}{1 - \frac{2\Phi}{c^2}}} \approx c \left(1 + \frac{2\Phi}{c^2}\right), \quad (3.1.3)$$

d'où l'indice de réfraction,

$$n = \frac{c}{c'} = \frac{1}{1 + \frac{2\Phi}{c^2}} \approx 1 - \frac{2\Phi}{c^2}. \quad (3.1.4)$$

En dérivant les extrema de (3.1.2) à partir des équations d'Euler Lagrange, on trouve la déviation totale de la trajectoire lumineuse causée par le déflecteur,

$$\hat{\alpha} = \frac{2}{c^2} \int_{l_A}^{l_B} \nabla_{\perp} \Phi dl. \quad (3.1.5)$$

Comme  $\Phi/c^2 \ll 1$ , les angles de déviation sont assez petits pour adopter l'approximation de Born, laquelle équivaut à intégrer sur le trajet non perturbé :

$$\hat{\alpha} = \frac{2}{c^2} \int_{-\infty}^{+\infty} \nabla_{\perp} \Phi dz, \quad (3.1.6)$$

où le chemin a été reparamétrisé par le décalage vers le rouge  $z$  par rapport au déflecteur.

Pour une masse ponctuelle,

$$\Phi = -\frac{GM}{r} \implies \nabla_{\perp} \Phi = \frac{GM}{r^3} \vec{x}_{\perp}, \quad (3.1.7)$$



où  $G$  est la constante gravitationnelle,  $M$  est la masse et  $r = \|\vec{x}\|$  est la distance de l'origine en 3D. L'angle de déviation prend donc la forme :

$$\hat{\alpha}(r) = \frac{4GM}{c^2 r_{\perp}^2} \vec{x}_{\perp}, \quad (3.1.8)$$

où  $r_{\perp} = \|\vec{x}_{\perp}\|$  est la distance de l'origine en 2D. Dans le cas d'un défecteur étendu, l'assertion que le défecteur est beaucoup plus petit que chaque section de la ligne de visée justifie l'approximation de l'« écran mince ». Celle-ci consiste à projeter sur des plans la distribution de masse du défecteur et le profil de lumineux de la source. La densité de surface du défecteur s'écrit donc :

$$\Sigma(\boldsymbol{\xi}) = \int \rho(\boldsymbol{\xi}, z) dz, \quad (3.1.9)$$

où  $\boldsymbol{\xi}$  est la position dans le plan du défecteur et  $\rho$  est la densité de masse tridimensionnelle. Les contributions de chaque élément de masse se somment pour donner l'angle de déviation total,

$$\hat{\alpha}(\boldsymbol{\xi}) = \frac{4G}{c^2} \int \frac{(\boldsymbol{\xi} - \boldsymbol{\xi}') \Sigma(\boldsymbol{\xi}')}{|\boldsymbol{\xi} - \boldsymbol{\xi}'|^2} d^2 \boldsymbol{\xi}', \quad (3.1.10)$$

### 3.1.2. Équation de lentille

La figure 3.1 schématise une lentille gravitationnelle.  $D_s$  désigne la distance angulaire de diamètre comobile entre le télescope et la source;  $D_d$ , celle entre le télescope et le défecteur; et  $D_{ds}$ , celle entre le défecteur et la source. Le trait plein illustre le trajet d'un rayon lumineux parti au point  $\boldsymbol{\eta}$  du plan de la source et dévié au point d'impact  $\boldsymbol{\xi}$  sur le plan du défecteur. La ligne formée de tirets montre le parcours d'un rayon lumineux en l'absence du défecteur. Le trait pointillé représente l'origine du système de coordonnées qui relie le télescope au centre du défecteur. La ligne pointillée entrecoupée de tirets indique la trajectoire apparente du rayon lumineux au-delà du défecteur.

Par la définition de la distance angulaire de diamètre,

$$\boldsymbol{\eta} = \frac{D_s}{D_d} \boldsymbol{\xi} - D_{ds} \hat{\alpha}(\boldsymbol{\xi}). \quad (3.1.11)$$

En exprimant les positions  $\boldsymbol{\eta}$  et  $\boldsymbol{\xi}$  par les coordonnées angulaires  $\boldsymbol{\beta}$  et  $\boldsymbol{\theta}$ , soit  $\boldsymbol{\eta} = D_s \boldsymbol{\beta}$  et  $\boldsymbol{\xi} = D_d \boldsymbol{\theta}$ , en plus de définir :

$$\boldsymbol{\alpha} \equiv \frac{D_{ds}}{D_s} \hat{\alpha}(D_d \boldsymbol{\theta}), \quad (3.1.12)$$

l'équation (3.1.11) devient :

$$\boldsymbol{\beta}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{D_{ds}}{D_s} \hat{\alpha}(D_d \boldsymbol{\theta}) = \boldsymbol{\theta} - \boldsymbol{\alpha}(\boldsymbol{\theta}). \quad (3.1.13)$$

L'équation (3.1.13) se nomme « équation de lentille » et elle dicte le tracé des rayons lumineux en présence d'un défecteur.

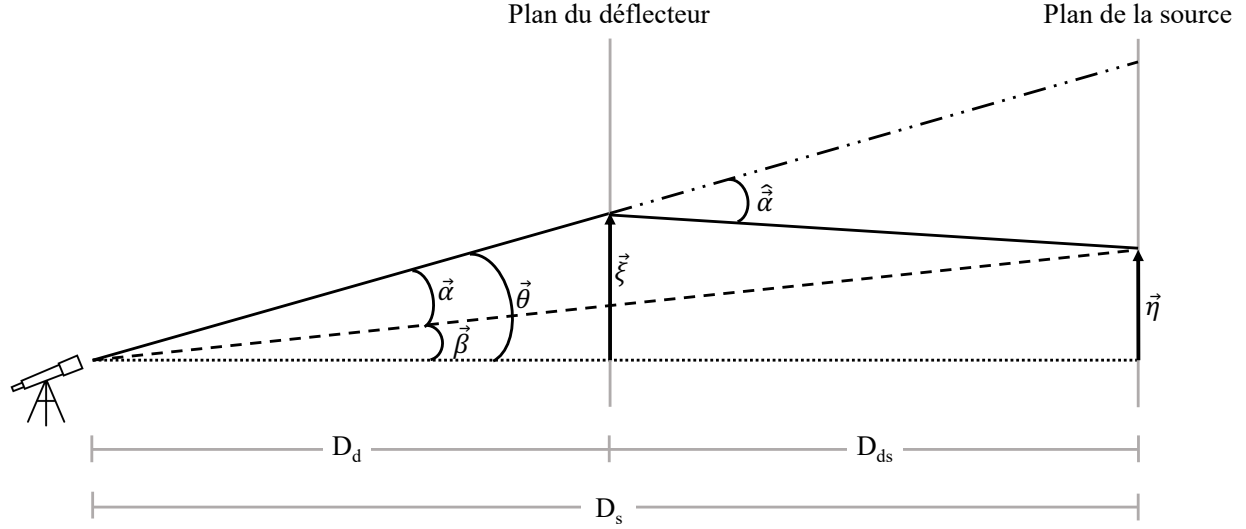


Fig. 3.1. Schéma du formalisme de lentille gravitationnelle

### 3.1.3. Potentiel de lentille

Un potentiel gravitationnel  $\Phi$  génère un potentiel de lentille  $\psi$  selon :

$$\psi(\boldsymbol{\theta}) = \frac{2D_{ds}}{D_d D_s c^2} \int \Phi(D_d \boldsymbol{\theta}) dz. \quad (3.1.14)$$

Le potentiel de lentille s'exprime donc aussi en fonction de la densité massique de surface adimensionnelle,

$$\kappa(\boldsymbol{\theta}) = \frac{\Sigma(D_d \boldsymbol{\theta})}{\Sigma_{cr}} \quad \text{où} \quad \Sigma_{cr} = \frac{c^2}{4\pi G} \frac{D_s}{D_d D_{ds}}, \quad (3.1.15)$$

par l'équation de Poisson,

$$\nabla^2 \psi(\boldsymbol{\theta}) = 2\kappa(\boldsymbol{\theta}). \quad (3.1.16)$$

Le potentiel de lentille détermine également les angles de déviation selon :

$$\boldsymbol{\nabla} \psi(\boldsymbol{\theta}) = \boldsymbol{\alpha}(\boldsymbol{\theta}). \quad (3.1.17)$$

La déformation des images et leur amplification découlent aussi du potentiel de lentille. La distorsion se décrit comme une matrice Jacobienne,

$$A(\boldsymbol{\theta}) = \frac{\partial \boldsymbol{\beta}}{\partial \boldsymbol{\theta}} = \delta_{ij} - \frac{\partial^2 \psi(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} = \begin{bmatrix} 1 - \kappa - \gamma_1 & -\gamma_2 \\ -\gamma_2 & 1 - \kappa + \gamma_1 \end{bmatrix}, \quad (3.1.18)$$

où le cisaillement  $\gamma = \gamma_1 + i\gamma_2$  est défini comme :

$$\gamma_1 = \frac{1}{2} \left( \frac{\partial^2 \psi(\boldsymbol{\theta})}{\partial^2 \theta_1} - \frac{\partial^2 \psi(\boldsymbol{\theta})}{\partial^2 \theta_2} \right); \quad \gamma_2 = \frac{\partial^2 \psi(\boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2}. \quad (3.1.19)$$

L'amplification s'obtient à partir de la matrice de distortion selon :

$$\mu(\boldsymbol{\theta}) = \frac{1}{\det A} = \frac{1}{(1 - \kappa)^2 - |\gamma|^2}. \quad (3.1.20)$$

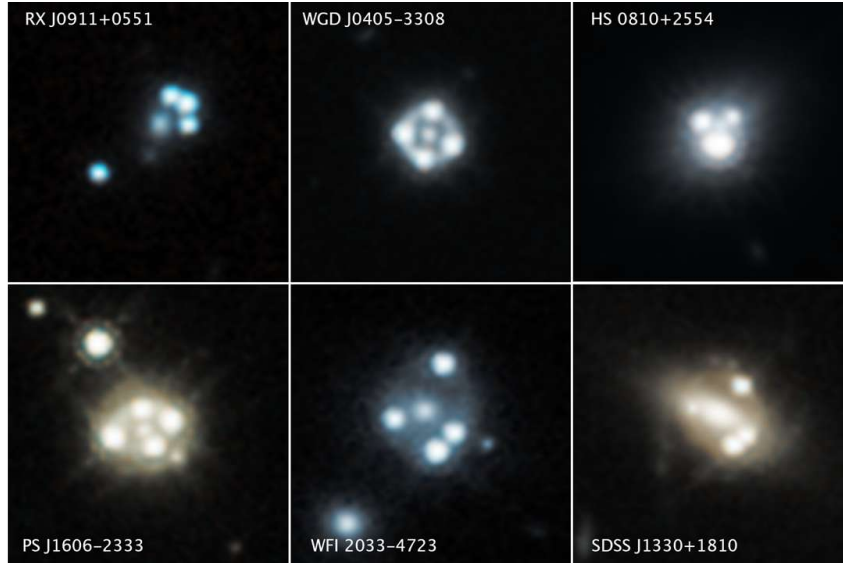
## 3.2. Applications

Les lentilles gravitationnelles sont de précieux outils de recherche en astrophysique. Arthur S. Eddington (1919) le prouve en montrant que la déviation de la lumière d'une étoile par le Soleil concorde davantage avec la prédiction de la relativité générale qu'avec celle de la mécanique classique. Chwolson (1924) et Einstein (1936) s'intéressent à l'effet de lentille gravitationnelle d'une étoile sur une autre et déduisent que l'image prend la forme d'un anneau lumineux. Einstein dérive aussi les positions des images et leur magnification. Malgré tout, il conclut qu'un tel système est pratiquement impossible à détecter. Zwicky (1937) considère plutôt des systèmes composés de deux galaxies. Il suggère de s'en servir pour mesurer les masses des galaxies ainsi que pour observer des galaxies lointaines dont les images sont petites et ténues.

### 3.2.1. Observer les sources lumineuses

Les lentilles gravitationnelles ont bel et bien amplifié des images de galaxies lointaines dont la luminosité n'aurait pas passé le seuil de détection autrement. Lorsque leurs images atteignent la Terre, ces galaxies apparaissent tel qu'elles étaient au moment où leur lumière a été émise. Ainsi, plus la source est éloignée, plus la lentille gravitationnelle révèle un passé lointain. La combinaison de ce phénomène aux instruments de pointe, comme le télescope spatial James Webb, donne donc accès aux premiers stades d'évolution des galaxies (Treu et al., 2022).

De plus, la magnification décuple la résolution spatiale des images, dévoilant des détails inédits de la morphologie des galaxies. Cette magnification est parfois si intense qu'elle permet de tester des modèles d'atmosphère stellaire à partir d'images d'étoiles lentillées. L'étude des galaxies hôtes d'un noyau actif bénéficie particulièrement de cet avantage. Les noyaux actifs sont des régions compactes et extrêmement lumineuses qui se forment autour d'un trou noir au centre d'une galaxie. La luminosité d'un noyau actif excède largement celle de sa galaxie hôte, ce qui dilue le signal de cette dernière. Or, sous la loupe d'une lentille gravitationnelle, la galaxie déformée en arcs ou en anneau se distingue facilement des images compactes du noyau actif qui est une source ponctuelle (voir Fig. 3.2). Il est alors possible d'analyser à la fois la galaxie hôte et son noyau actif avec une meilleure résolution, surtout si le noyau actif subit aussi un effet de microlentille. (Schneider, 2006)



**Fig. 3.2.** Six lentilles gravitationnelles fortes capturées par le télescope spatial Hubble impliquant des noyaux actifs et leur galaxie hôte. Crédit : NASA, ESA, A. Nierenberg (JPL) and T. Treu (UCLA)

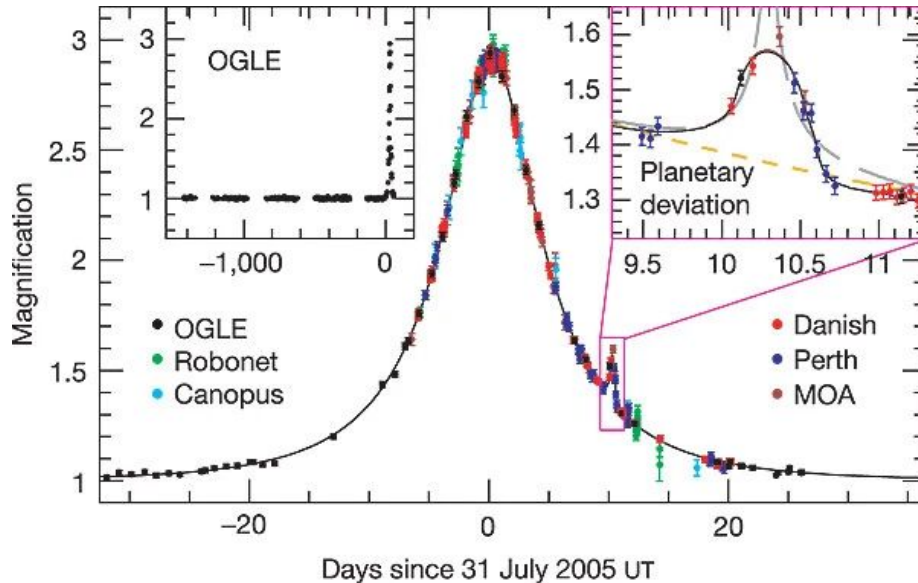
### 3.2.2. Étudier les défecteurs massifs

Comme le prévoyait Zwicky, les lentilles gravitationnelles ont effectivement fourni des mesures précises de masses galactiques. En fait, seuls les effets gravitationnels rendent compte de la masse totale des galaxies, puisque celles-ci se composent majoritairement de matière sombre. Cette substance compose 84% de toute l'énergie sous forme de matière dans l'Univers (Collaboration Planck et al., 2020). Or, sa nature reste inconnue, notamment parce qu'elle n'interagit pas électromagnétiquement, ce qui la rend indétectable par des moyens optiques. En revanche, toute matière qui possède une masse contribue à l'effet de lentille gravitationnelle, d'où la pertinence de ce phénomène pour évaluer la masse des galaxies (Massey et al., 2010).

L'effet de lentille gravitationnelle, dans ses versions milli et micro, trahit également la présence de fines structures au sein du défecteur et le long de la ligne de visée. Par exemple, le potentiel de lentille d'une galaxie se décrit adéquatement par un modèle lisse, mais les sous-halos de matière sombre y introduisent des perturbations. L'image lentillée renseigne donc sur la distribution des sous-halos et leur nombre, ce qui permet d'étudier le problème des satellites manquants. (Treu, 2010)

Les planètes agissent également comme des microlentilles. Dans le cas où deux étoiles s'alignent sur la ligne de visée d'un télescope, celle d'avant-plan produit un effet de lentille sur l'image de celle d'arrière-plan. Si un membre du système planétaire de la première se trouve au bon endroit au bon moment, il contribue à l'effet de lentille. Un pic de brillance

secondaire, souvent subtil et de courte durée, surgit sur celui généré par l'étoile, permettant ainsi de détecter des exoplanètes (voir Fig. 3.3; Mao & Paczynski, 1991; Gould & Loeb, 1992; Bond et al., 2004). Selon les archives de la NASA <sup>1</sup>, l'effet de microlentille a permis de détecter 176 exoplanètes jusqu'à maintenant.



**Fig. 3.3.** Détection d'une planète par son effet de microlentille. Les points indiquent la magnification d'une étoile par l'effet de lentille d'une autre, telle que mesurée par six relevés astronomiques à différents moments. Le graphique en haut à droite agrandit la période de magnification supplémentaire causée par une planète. La ligne pleine noire représente le meilleur modèle de déflecteur binaire composé d'une étoile et d'une planète. La ligne pointillée grise désigne le meilleur modèle de source binaire (rejeté par les données). La ligne pointillée orange montre le meilleur modèle de déflecteur unique, c'est-à-dire une étoile sans planète. Crédit : Beaulieu et al. (2006)

### 3.2.3. Inférer la constante de Hubble

Refsdal (1964) propose une approche basée sur les lentilles gravitationnelles fortes pour mesurer le taux d'expansion de l'Univers. Il se concentre sur un système ayant une supernova comme source lumineuse et une galaxie comme déflecteur. Il assume aussi que le déflecteur et la source s'alignent de façon à ce que la gravitation du premier génère deux images de la seconde. La supernova étant une source ponctuelle, les deux images sont clairement séparées. De plus, Refsdal considère que la supernova émet un pic de luminosité, lequel se transfère aux deux images. Toutefois, les parcours des deux images ont des longueurs, et donc des durées,

<sup>1</sup>[https://exoplanetarchive.ipac.caltech.edu/docs/counts\\_detail.html](https://exoplanetarchive.ipac.caltech.edu/docs/counts_detail.html)

différentes. Le pic lumineux apparaît donc à des moments distincts sur chaque image. Le délai entre l'arrivée du pic à l'image A et son arrivée à l'image B s'exprime comme :

$$\Delta t = \frac{D_{\Delta t}}{c} [\phi(\boldsymbol{\theta}_A, \boldsymbol{\beta}) - \phi(\boldsymbol{\theta}_B, \boldsymbol{\beta})], \quad (3.2.1)$$

où  $D_{\Delta t}$  est la *time-delay distance* et  $\phi$  est le potentiel de Fermat. Ce dernier est déterminé seulement par le système de lentille gravitationnelle et il est défini comme :

$$\phi(\boldsymbol{\theta}, \boldsymbol{\beta}) \equiv \frac{(\boldsymbol{\theta} - \boldsymbol{\beta})^2}{2} - \psi(\boldsymbol{\theta}). \quad (3.2.2)$$

Quant à la *time delay distance*, elle s'écrit :

$$D_{\Delta t} \equiv (1 + z_d) \frac{D_d D_s}{D_{ds}}, \quad (3.2.3)$$

où  $z_d$  est le décalage vers le rouge du déflecteur. En assumant une courbure spatio-temporelle nulle, la distance de diamètre angulaire entre deux décalages vers le rouge  $z_1$  et  $z_2$  se calcule selon :

$$D(z_1, z_2) = \frac{1}{1 + z_2} \int_{z_1}^{z_2} dz' \frac{c}{H_0 \sqrt{\Omega_m (1 + z')^3 + \Omega_\Lambda}}, \quad (3.2.4)$$

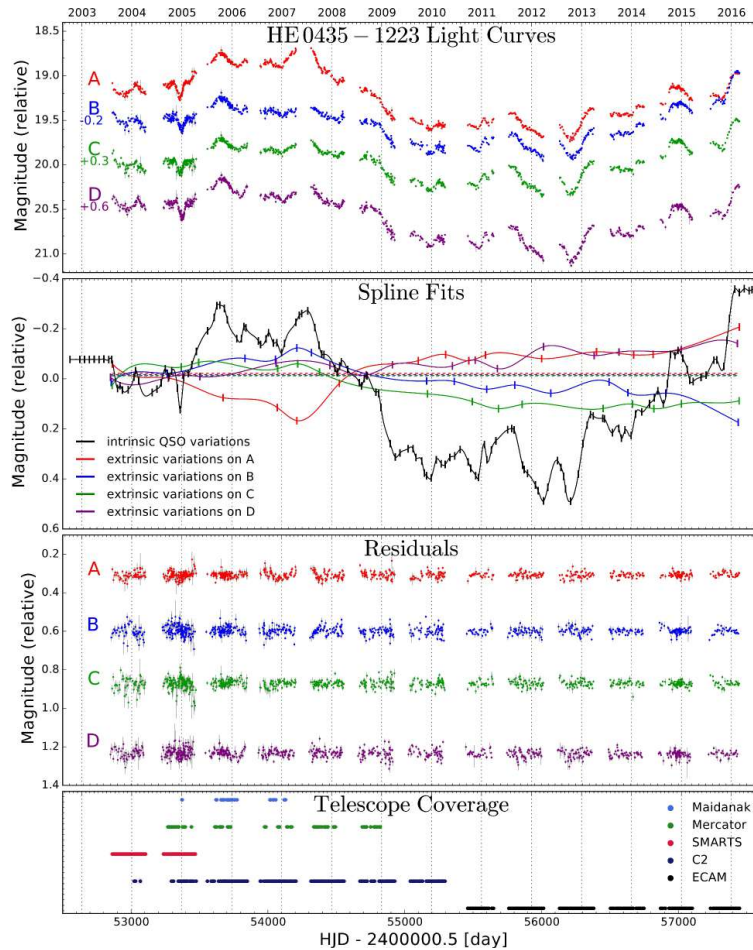
avec  $\Omega_m$  pour la densité d'énergie de la matière et  $\Omega_\Lambda$  pour celle de l'énergie sombre.  $D_{\Delta t}$  dépend donc étroitement de  $H_0$ , mais faiblement des autres paramètres cosmologiques.

La mise en application de cette méthode devient possible alors que Walsh et al. (1979), profitant des toutes nouvelles CCD et du *Very Large Array* (VLA), observent pour la première fois deux images d'un quasar fortement lentillé par une galaxie d'avant-plan. Le formalisme développé ci-haut demeure valide lorsqu'un quasar remplace la supernova. En effet, les quasars, qui sont des noyaux actifs de galaxies, se substituent aux supernovae puisque, tout comme elles, ils sont compacts et leur luminosité est variable. Toutefois, l'inférence de la constante de Hubble par ce procédé compte encore de nombreux défis.

Pour inférer la constante de Hubble à partir d'une source variable fortement lentillée, il faut obtenir les décalages vers le rouge du déflecteur  $z_d$  et de la source  $z_s$  par spectroscopie. Un suivi de la courbe de lumière de chaque image doit être effectué afin de les comparer, et ainsi, mesurer les délais  $\Delta t$  (eq. 3.2.1). Une modélisation de la densité de masse du déflecteur  $\kappa$  (eq. 3.1.15) est nécessaire pour déterminer la différence de potentiel de Fermat  $\phi$  (eq. 3.2.2) entre les images. Finalement, briser les dégénérescences du modèle de masse nécessite la vitesse de dispersion du déflecteur et des observables dans son environnement.

Au niveau des courbes de lumière, l'effet de microlentille pose problème s'il est puissant ou si sa fréquence s'approche de celle de variation du quasar. Il enterre alors le signal du quasar ou se confond avec lui. Une cadence d'observation soutenue ainsi qu'un ratio élevé entre le signal et le bruit peuvent capturer les fluctuations du quasar à des échelles temporelles inférieures à celles induites par l'effet de microlentille. Les deux variables se distingueraient alors

davantage. Remplacer les quasars par des supernovae aiderait également, puisque la courbe de lumière de ces dernières est régulière, et donc facilement identifiable, contrairement à celle des quasars qui est aléatoire.



**Fig. 3.4.** Le graphique du haut montre les courbes de lumière des quatre images du quasar HE 0435-1223. Le second graphique présente un modèle de la variation intrinsèque du quasar en noir et les courbes de variation extrinsèque de chaque image. Le troisième graphique affiche les valeurs résiduelles. Sur le graphique du bas, chaque point indique une nuit d'observation, et ce, pour les cinq télescopes impliqués durant les 13 années de suivi. Crédit : Bonvin et al. (2016)

De plus, en raison des saisons, un objet astronomique se trouve dans le champ de vision d'un observatoire terrestre entre 6 et 9 mois par année selon sa position dans le ciel. Cette restriction introduit des discontinuités dans les courbes de lumière enregistrées. À ces difficultés s'ajoutent un ratio variable entre le signal et le bruit, en plus d'effets atmosphériques. Avec la qualité actuelle des données, une corrélation des courbes de lumière par translation temporelle donne tout de même des mesures exactes et précises des délais, mais cette stratégie exige parfois une décennie de suivi (Suyu et al., 2018). La figure 3.4 illustre ces complications

en présentant les courbes de lumière des quatre images du quasar HE 0435-1223.

En ce qui concerne la modélisation de la distribution de masse causant l'effet de lentille, plusieurs limitations subsistent. Assumer un modèle paramétrique simple fonctionne si la résolution de l'image est faible, mais ne suffit pas lorsque les effets des sous-structures du déflecteur apparaissent distinctement. De plus, le procédé consiste à simuler l'observation selon des paramètres initiaux, puis à les optimiser pour maximiser une fonction de vraisemblance. Cette analyse requiert beaucoup de temps et doit s'effectuer individuellement sur chaque observation. Outre les modèles paramétriques, il existe aussi des représentations pixelisées, mais définir une distribution a priori sur ces paramétrisations est difficile, ce qui les limite par rapport à leurs homologues continues.

Une difficulté supplémentaire provient de la dégénérescence de la « feuille de masse » (Falco et al., 1985; Schneider & Sluse, 2013, 2014; Schneider, 2019). Considérons un modèle de distribution de masse  $\kappa$  qui reproduit adéquatement une observation de lentille gravitationnelle, incluant la position des images, leurs flux relatifs et leurs formes. Ce  $\kappa$  appartient à une famille de modèles  $\kappa_\lambda$  qui satisfont :

$$\kappa_\lambda = \lambda + (1 - \lambda)\kappa, \quad (3.2.5)$$

où  $\lambda$  est une constante. Le premier terme correspond à l'ajout d'une densité surfacique uniforme au déflecteur, tandis que le second équivaut à redimensionner celle d'origine. Sous cette transformation, les angles de déviation deviennent :

$$\alpha_\lambda = (1 - \lambda)\theta + \lambda\alpha(\theta), \quad (3.2.6)$$

ce qui affecte l'équation de lentille selon

$$\beta_\lambda = \theta - \alpha_\lambda(\theta) \quad (3.2.7)$$

$$= \theta - [(1 - \lambda)\theta + \lambda\alpha(\theta)] \quad (3.2.8)$$

$$= \lambda\theta - \lambda\alpha(\theta) \quad (3.2.9)$$

$$= \lambda\beta \quad (3.2.10)$$

L'inférence de la position de la source s'effectue simultanément avec celle du modèle de masse. De ce fait, n'importe quel modèle de la famille décrite par l'équation (3.2.5), combiné à l'ajustement approprié des coordonnées de la source, concorde tout aussi bien avec l'observation. De plus, le cisaillement réduit,

$$g = \frac{\gamma}{1 - \kappa}, \quad (3.2.11)$$

qui détermine la distorsion des images, reste invariant sous la transformation de la feuille de masse. Les modèles  $\kappa_\lambda$  ne se distinguent donc pas par la forme des images. Qui plus est,



bien que l’amplification se modifie comme :

$$\mu_\lambda = \frac{\mu}{\lambda^2}, \quad (3.2.12)$$

les flux relatifs demeurent les mêmes, ce qui ne permet guère plus de discriminer les  $\kappa_\lambda$ . En revanche, les potentiels de Fermat relatifs sont affectés selon :

$$\Delta\phi_\lambda = (1 - \lambda)\Delta\phi. \quad (3.2.13)$$

L’effet sur les délais entre les images lentillées d’un quasar se confond avec celui de la cosmologie, ce qui compromet l’inférence de  $H_0$  par cette méthode. Toutefois, la vitesse de dispersion stellaire au sein du défecteur brise la dégénérescence entre les  $\kappa_\lambda$ , puisqu’elle est sensible à la valeur intégrée de  $\kappa$  à l’intérieur du rayon effectif du défecteur. Aussi, une chandelle standard comme source, par exemple une supernova, restreindrait la valeur de  $\mu$ .

La dégénérescence de la feuille de masse émerge aussi des structures le long de la ligne de visée. Si leurs effets sont suffisamment petits, ces structures peuvent être traitées comme une projection de leur densité de masse surfacique dans le plan de la lentille. Le terme  $\kappa_{\text{ext}}$  regroupe leurs contributions et influence la *time-delay distance* selon :

$$D_{\Delta t} = \frac{D_{\Delta t}^{\text{model}}}{(1 - \kappa_{\text{ext}})}. \quad (3.2.14)$$

Une technique pour briser cette dégénérescence consiste à estimer statistiquement la densité de masse le long de la ligne de visée.

Malgré tous ces obstacles, les lentilles gravitationnelles constituent une option pertinente pour investiguer la crise de la cosmologie, puisque cette approche est indépendante du FDC et de l’échelle des distances cosmiques. Par ailleurs, le grand échantillon statistique que fourniront les prochains relevés astronomiques améliorera considérablement la précision de cette méthode. Oguri & Marshall (2010) prévoient 8000 observations de quasars lentillés pendant les 10 ans d’activité du LSST. Goldstein & Nugent (2016) prédisent 500 détections de supernovae lentillées pour la même campagne. De plus, l’apprentissage automatique fournit désormais de nombreuses techniques d’analyse plus rapides et plus robustes que les méthodes traditionnelles (Hezaveh et al., 2017; Levasseur et al., 2017; Morningstar et al., 2019; Pearson et al., 2019; Schuldt, S. et al., 2021; Park et al., 2021; Legin et al., 2021, 2022; Adam et al., 2022).



# Chapitre 4

---

## Apprentissage automatique

L'apprentissage automatique s'inscrit dans la branche appliquée des statistiques. Il s'appuie sur un vaste ensemble de données, sur un modèle flexible et sur un algorithme d'optimisation efficace pour construire une fonction qui exécute une tâche. Un réseau de neurones artificiels se prête bien à la représentation de la fonction. Il comprend un nombre fixe de fonctions paramétriques de base, mais laisse leurs paramètres s'adapter au cours de l'apprentissage. Parmi les tâches les plus couramment effectuées par des réseaux de neurones artificiels, on compte la régression, la classification, la traduction, la détection d'anomalie, l'échantillonnage, la segmentation d'images, le remplissage d'images, l'estimation de distributions, etc. L'introduction à l'apprentissage automatique suivante se base sur le livre de référence de Bishop (2006).

### 4.1. Perceptron multicouche

Un perceptron multicouche est un réseau de neurones artificiels constitué de plusieurs couches de neurones, soit une d'entrée, un nombre prédéterminé de cachées et une de sortie. La couche d'entrée représente les données initiales. Les couches cachées résultent de la combinaison linéaire des valeurs transmises par la couche précédente et d'une opération non-linéaire. Habituellement, le même nombre de neurones compose chaque couche cachée, ce nombre étant choisi arbitrairement ou par optimisation. La couche de sortie fournit les résultats finaux. Sa dimension dépend du nombre de données de sortie souhaité.

Considérons un vecteur de données d'entrée  $\mathbf{x}$  de dimension  $D$ , soit  $[x_1, x_2, \dots, x_D]$ , et une première couche cachée de  $H$  neurones. Une activation  $a_j$  se calcule pour chaque neurone  $j$  de cette couche selon :

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}, \quad (4.1.1)$$

où  $w_{ji}^{(1)}$ , nommé « poids », est un coefficient et où  $w_{j0}^{(1)}$  est appelé « biais ». Les indices  $i = 1, 2, \dots, D$  et  $j = 1, 2, \dots, H$  réfèrent à une dimension d'entrée et à un neurone de la couche

cachée, respectivement. L'exposant (1) indique l'appartenance à la première couche. Le réseau de neurones évalue ensuite une fonction d'activation  $h(\cdot)$ , non-linéaire et différentiable, pour chaque neurone, ce qui donne l'unité cachée,

$$z_j = h(a_j). \quad (4.1.2)$$

Des exemples de fonction d'activation sont tanh, sigmoïde, ReLU, Leaky ReLU et ELU (voir Fig. 4.1), respectivement définies comme

$$h(a_j) = \tanh(a_j) \quad (4.1.3)$$

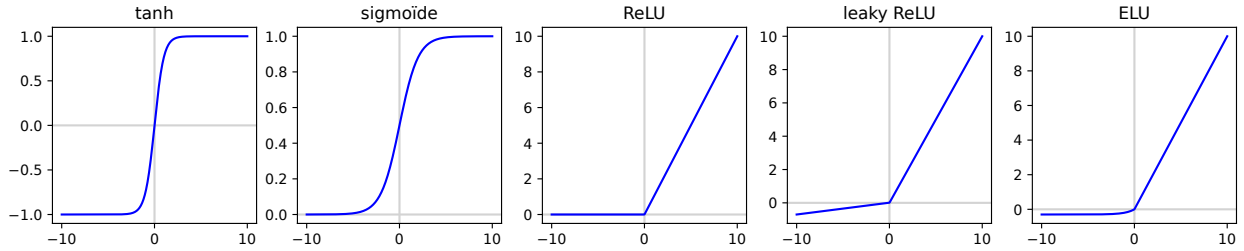
$$h(a_j) = \frac{1}{1 + e^{-a_j}} \quad (4.1.4)$$

$$h(a_j) = \max(0, a_j) \quad (4.1.5)$$

$$h(a_j) = \max(\epsilon a_j, a_j) \quad (4.1.6)$$

$$h(a_j) = \begin{cases} a_j & \text{pour } a_j \geq 0 \\ \epsilon(e^{a_j} - 1) & \text{pour } a_j < 0 \end{cases} \quad (4.1.7)$$

où  $0 < \epsilon < 1$  est un paramètre à ajuster. Les couches cachées ont communément recours aux fonctions ReLU, Leaky ReLU ou ELU. Pour la couche de sortie, le choix de fonction d'activation dépend de la nature de la valeur désirée. Les fonctions sigmoïde et tanh se prêtent bien à la prédiction de probabilités car leur image est confinée entre deux asymptotes. ReLU, Leaky ReLU et ELU servent davantage à la régression.



**Fig. 4.1.** Représentation des fonctions d'activation définies par les équations de (4.1.7) à (4.1.3).

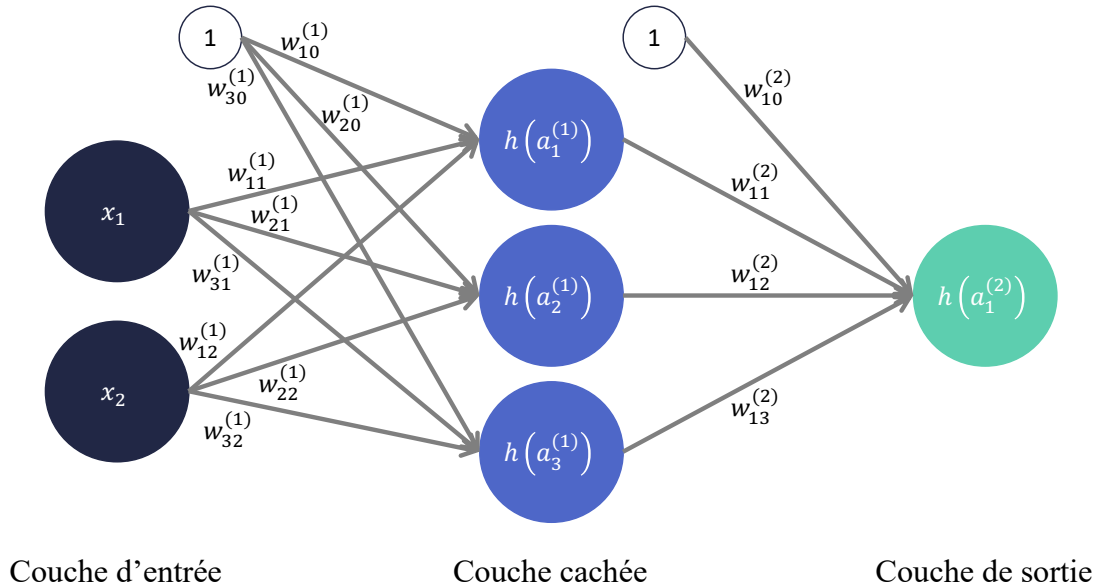
La figure 4.2 schématise le flot d'information dans un perceptron multicouche contenant une seule couche cachée. Sa sortie  $\hat{\mathbf{y}}$  de dimension  $K$  s'obtient par la fonction représentée par le réseau de neurones en entier,

$$\hat{y}_k(\mathbf{x}, \mathbf{w}) = h_s \left( \sum_{j=i}^M w_{kj}^{(2)} h_c \left( \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right), \quad (4.1.8)$$

où  $k = 1, \dots, K$  désigne une dimension de la donnée de sortie;  $\mathbf{w}$  représente tous les poids du réseau;  $h_s$  est la fonction d'activation suivant la couche de sortie et  $h_c$ , celle suivant la couche

cachée; et l'exposant (2) indique la seconde couche. Comme mentionné précédemment, les couches cachées appliquent généralement la même fonction d'activation, mais la couche de sortie peut en utiliser une autre.

Hornik (1991) a démontré qu'un perceptron multicouche est un approximateur universel. En d'autres mots, s'il a assez de neurones, il peut modéliser n'importe quelle fonction lisse aussi fidèlement que voulu. Un élément clé de son succès est qu'il est différentiable par rapport à ses poids  $w$ , ce qui permet l'ajustement de ceux-ci par entraînement.



**Fig. 4.2.** Schéma d'un perceptron multicouche dont l'entrée a deux dimensions et la sortie en a une. Il comprend une seule couche cachée de trois neurones. L'entrée unitaire dont les poids sont en fait les biais permet de simplifier la notation, voir équation (4.2.5).

## 4.2. Entraînement

Supposons qu'un perceptron multicouche est entraîné à évaluer le prix de vente d'une maison à partir de son âge, de son nombre de pièces et de la superficie de son terrain. L'ensemble de données disponible comprend les caractéristiques et le prix de plusieurs maisons récemment vendues. L'entraînement sera alors supervisé, ce qui signifie que le réseau de neurones artificiels apprend à partir d'exemples de données d'entrée et de leurs cibles associées, ces dernières étant les valeurs à prédire.

L'entraînement se déroule comme suit. Les caractéristiques des maisons sont introduites en entrée du réseau de neurones, lequel propage l'information jusqu'à sa couche de sortie pour faire ses prédictions. Une fonction de perte compare ensuite les estimations du perceptron multicouche avec les véritables prix correspondants. Les gradients de la fonction de perte

par rapport aux poids du réseau de neurones sont calculés par rétropropagation. Une descente de gradient effectue alors un pas dans l'espace des poids afin de les optimiser. Ces étapes se répètent jusqu'à convergence. Une fois l'entraînement terminé, la performance du perceptron multicouche s'évalue sur des exemples qui ne faisaient pas partie de ses exemples d'entraînement.

### 4.2.1. Fonction de perte

L'objectif de l'entraînement est de minimiser la fonction de perte, laquelle traduit mathématiquement la tâche à accomplir. Dans la plupart des cas, elle provient d'une assertion sur la distribution des cibles, soit la « fonction de vraisemblance ». Cette dernière quantifie la probabilité d'une observation étant donné un modèle et ses paramètres. Maximiser une fonction de vraisemblance équivaut à trouver le modèle et ses paramètres pour lesquels l'observation est la plus probable. Il est pratique de prendre le logarithme de la fonction de vraisemblance, puisqu'il a le même maximum qu'elle, mais qu'il est moins abrupte. De plus, puisque la convention est de minimiser plutôt que de maximiser, on se sert du négatif du logarithme comme fonction de perte.

Prenons l'exemple d'une classification binaire. Dans ce cas, les cibles sont des étiquettes  $y = 1$  et  $y = 0$  qui désignent respectivement les classes  $\mathcal{C}_1$  et  $\mathcal{C}_2$ . Une prédiction du réseau de neurones s'interprète alors comme la probabilité qu'un exemple donné appartienne à l'une des classes :

$$p(\mathcal{C}_1 | \mathbf{x}) = \hat{y}(\mathbf{x}, \mathbf{w}) \quad (4.2.1)$$

$$p(\mathcal{C}_2 | \mathbf{x}) = 1 - \hat{y}(\mathbf{x}, \mathbf{w}) . \quad (4.2.2)$$

La probabilité des cibles étant donné les données d'entrée prend donc la forme d'une distribution de Bernoulli,

$$p(y | \mathbf{x}, \mathbf{w}) = \hat{y}(\mathbf{x}, \mathbf{w})^y [1 - \hat{y}(\mathbf{x}, \mathbf{w})]^{1-y} . \quad (4.2.3)$$

Assumant des exemples indépendants et identiquement distribués (iid), la fonction de perte obtenue à partir de cette fonction de vraisemblance est l'entropie croisée. Pour un ensemble de  $N$  exemples, celle-ci s'écrit :

$$\mathcal{L}(\mathbf{w}) = - \sum_{n=1}^N [y_n \ln \hat{y}_n(\mathbf{x}, \mathbf{w}) + (1 - y_n) \ln(1 - \hat{y}_n(\mathbf{x}, \mathbf{w}))] , \quad (4.2.4)$$

où  $n = 1, \dots, N$  identifie les exemples d'entraînement.

Une fois que la fonction de perte est évaluée, il faut calculer son gradient par rapport à tous les poids du réseau de neurones. Une procédure efficace consiste à propager les gradients de la couche de sortie vers celle d'entrée.

## 4.2.2. Rétropropagation des gradients

L'information se propage dans un réseau de neurones artificiels en calculant les activations et les unités cachées successivement de la couche d'entrée à la couche de sortie. Alors, au moment de calculer les gradients de la fonction de perte par rapport aux poids du réseau de neurones, toutes les activations  $a_j$  et toutes les unités cachées  $z_i$  sont des valeurs connues. Commençons par réécrire l'équation (4.1.1) pour les couches cachées du réseau de neurones,

$$a_j = \sum_i w_{ji} z_i, \quad (4.2.5)$$

où le biais  $w_{j0}$  a été absorbé par  $w_{ji}$  et  $z_i$  a un élément unitaire de plus, ce qui équivaut à l'addition explicite du biais de l'équation (4.1.1).

Ensuite, en assumant des exemples d'entraînement iid, la fonction de perte totale s'écrit comme la somme de ses évaluations sur chacun des exemples individuels :

$$\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \mathcal{L}_n(\mathbf{w}). \quad (4.2.6)$$

Ainsi, le gradient total est aussi la somme des gradients sur chacun des exemples. L'un d'entre eux s'écrit :

$$\frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial w_{ji}} = \frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}}. \quad (4.2.7)$$

Le second terme de droite s'obtient par l'équation (4.2.5),

$$\frac{\partial a_j}{\partial w_{ji}} = z_i, \quad (4.2.8)$$

pour un  $i$  particulier. En insérant (4.2.8) dans (4.2.7), le gradient devient :

$$\frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial w_{ji}} = \frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_j} z_i. \quad (4.2.9)$$

Le premier terme de droite des équations (4.2.7) et (4.2.9) s'exprime par une dérivée en chaîne,

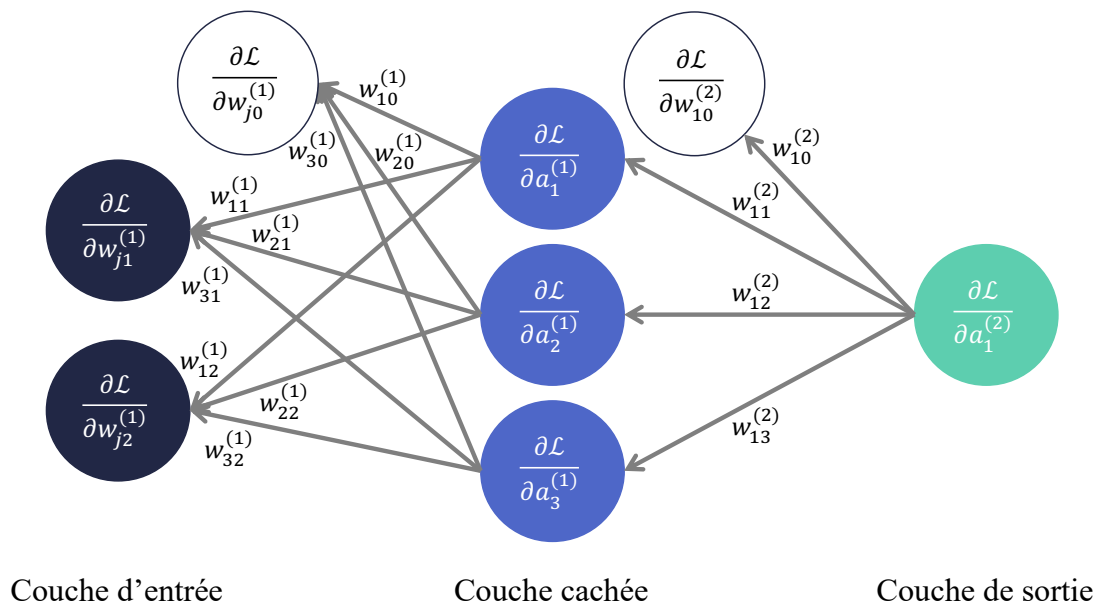
$$\frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_j} = \sum_k \frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_k} \frac{\partial a_k}{\partial a_j}. \quad (4.2.10)$$

En se référant aux équations (4.2.8) et (4.1.2), on trouve :

$$\frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_j} = \frac{\partial h(a_j)}{\partial a_j} \sum_k w_{kj} \frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_k}. \quad (4.2.11)$$

Ainsi, au moment d'évaluer les gradients, on calcule d'abord  $\frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_j}$  pour la couche de sortie. Cette information sert alors à calculer la même quantité, mais pour la couche précédente, et ainsi de suite. L'information se « rétropropage » donc selon l'équation (4.2.11) à travers l'ensemble des couches du réseau. Toutes les valeurs de  $\frac{\partial \mathcal{L}_n(\mathbf{w})}{\partial a_j}$  se déterminent de la sorte.

Puis, les gradients sont finalement obtenus grâce à l'équation (4.2.9). La figure 4.3 illustre ce processus. Pour terminer une phase d'entraînement, il ne reste qu'à actualiser les poids.



**Fig. 4.3.** Schéma de la rétropropagation du gradient.

### 4.2.3. Descente de gradient

L'entraînement commence par l'initialisation des poids. Pour ce faire, Glorot & Bengio (2010) proposent de les tirer de la distribution,

$$w_{ji} \sim \mathcal{U} \left( -\sqrt{\frac{6}{D+K}}, \sqrt{\frac{6}{D+K}} \right), \quad (4.2.12)$$

où  $D$  et  $K$  sont les dimensions d'entrée et de sortie, respectivement.

Le réseau de neurones procède ensuite à sa première prédiction, sur laquelle s'évalue la fonction de perte. Cette dernière se décompose en une somme sur les exemples, en supposant que ceux-ci sont iid. Ainsi, la fonction de perte équivaut à une espérance mathématique sur la distribution empirique définie par l'ensemble d'entraînement. Son gradient exact se calcule donc sur l'entièreté des données.

En pratique, l'ensemble d'entraînement est plutôt divisé en lots, soit des sous-ensembles de données de petite taille et tirés aléatoirement d'une distribution uniforme. Une étape d'entraînement, appelée « époque », consiste à estimer le gradient et à effectuer un pas dans l'espace des paramètres séquentiellement à partir de chaque lot.

De cette façon, l'algorithme converge en moins d'opérations, puisque se fier à une approximation du gradient pour effectuer les pas ne le désavantage pas autant que le nombre de



calculs requis par le gradient exact. Cette stratégie économise aussi des opérations si les exemples qui ont une contribution similaire au gradient sont séparés dans différents lots. De plus, elle sort plus facilement des minima locaux, puisque ceux-ci changent d'un lot à l'autre, ce qui améliore la capacité du réseau de neurones à généraliser.

Pour un lot  $\mathcal{B} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(B)}, y^{(B)})\}$ , le gradient s'écrit :

$$\mathbf{g} = \frac{1}{B} \nabla_{\mathbf{w}} \sum_{b=1}^B \mathcal{L}(\mathbf{x}^{(b)}, y^{(b)}, \mathbf{w}). \quad (4.2.13)$$

Les algorithmes d'optimisation les plus fréquemment utilisés en apprentissage automatique sont la descente de gradient stochastique et ses variantes. Dans sa version de base, les poids du réseau de neurones sont mis à jour selon :

$$\mathbf{w}_{\tau} = \mathbf{w}_{\tau-1} - \eta_{\tau} \mathbf{g}, \quad (4.2.14)$$

où  $\tau$  indique l'itération actuelle et où  $\eta_{\tau}$  est le taux d'apprentissage. Ce dernier sert à diminuer la grandeur du pas dans l'espace des poids afin de ne pas dépasser le minimum. Pour cette raison, il est recommandé de le diminuer tandis que l'entraînement progresse, c'est-à-dire plus les pas s'approchent du point optimal, d'où l'indice  $\tau$ .

Une astuce pour accélérer l'apprentissage, proposée par Polyak (1964), est d'avoir recours au momentum. Il s'agit d'une analogie avec la mécanique classique où l'hypersurface de la fonction de perte dans l'espace des poids se compare à un potentiel gravitationnel dans lequel se déplace une particule. Cette technique consiste donc à orienter le pas en fonction d'une moyenne mobile exponentiellement décroissante des gradients précédents. Pour ce faire, on introduit une vitesse  $\mathbf{v}$ , ou momentum,

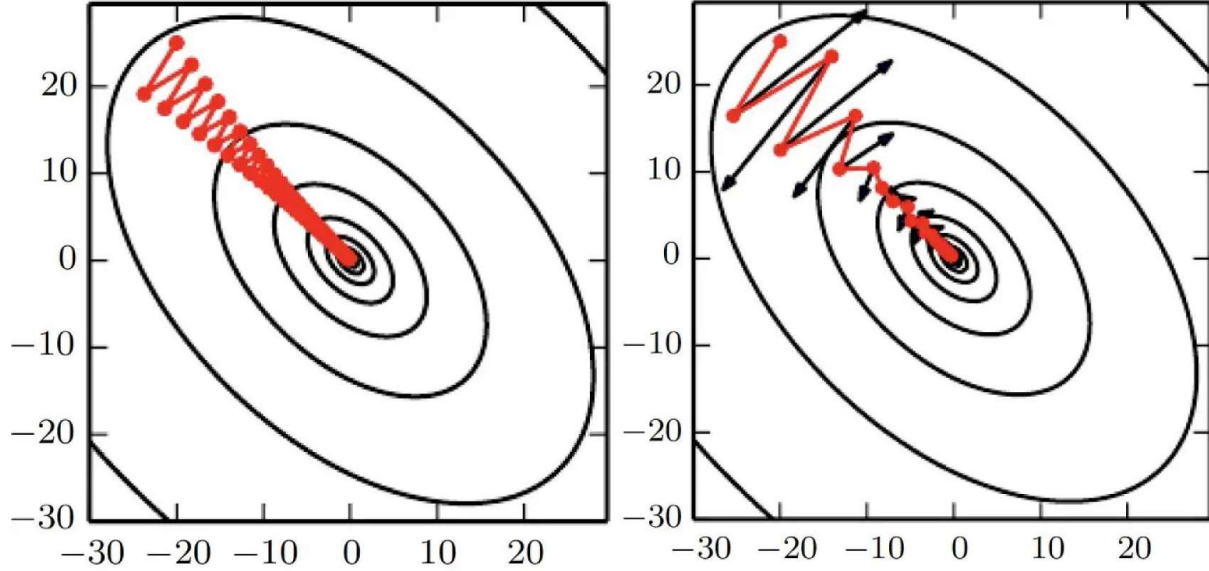
$$\mathbf{v}_{\tau} = \omega \mathbf{v}_{\tau-1} - \eta_{\tau} \mathbf{g}, \quad (4.2.15)$$

où  $\omega$  est le taux de décroissance. Puis, les poids s'actualisent selon :

$$\mathbf{w}_{\tau} = \mathbf{w}_{\tau-1} + \mathbf{v}_{\tau}. \quad (4.2.16)$$

Le pas sera donc d'autant plus grand que les gradients précédents s'alignent, comme une balle gagne en vitesse en dévalant une pente. La figure 4.4 illustre cet effet.

D'autres algorithmes améliorent l'efficacité de l'entraînement en introduisant plutôt un taux d'apprentissage adaptatif. Duchi et al. (2011) conçoivent ADAGRAD, lequel assigne un taux d'apprentissage spécifique à chaque poids et le réduit d'un facteur inversement proportionnel à la norme de tous les gradients passés. Or l'utilisation de l'historique complet des gradients risque de diminuer les pas trop rapidement et trop intensément, l'empêchant de rejoindre un minimum. Hinton et al. (2012) présentent plutôt RMSPROP. Tout comme ADAGRAD, la



**Fig. 4.4.** Représentation de l'effet du momentum. Les ellipses noires concentriques indiquent les isocontours d'une fonction quadratique bidimensionnelle. À gauche, les points reliés en rouge suivent la trajectoire d'une descente de gradient stochastique et, à droite, celle d'une descente de gradient avec du momentum. Les flèches noires montrent le pas qu'une descente stochastique aurait fait à ce point. La descente de gradient avec du momentum requiert beaucoup moins de points pour rejoindre le minimum. Crédit : Goodfellow et al. (2016)

mise à jour des poids s'écrit :

$$\mathbf{w}_\tau = \mathbf{w}_{\tau-1} + \frac{\eta_\tau}{\delta + \sqrt{\mathbf{r}_\tau}}, \quad (4.2.17)$$

où  $\mathbf{r}_\tau$  est l'accumulation des gradients et  $\delta$  est une faible valeur à des fins de stabilité numérique. L'amélioration apportée à ADAGRAD par RMSPROP est l'emploi d'une moyenne pondérée, mobile et exponentiellement décroissante des gradients plutôt que leur norme cumulative. En termes mathématiques, ADAGRAD accumule les gradients comme :

$$\mathbf{r}_\tau = \mathbf{r}_{\tau-1} + \mathbf{g} \odot \mathbf{g}, \quad (4.2.18)$$

tandis que RMSPROP fait de même selon :

$$\mathbf{r}_\tau = \omega \mathbf{r}_{\tau-1} + (1 - \omega) \mathbf{g} \odot \mathbf{g}, \quad (4.2.19)$$

où  $\omega$  est encore un taux de décroissance. De cette façon, l'impact des gradients les plus anciens disparaît progressivement.

ADAMAX, développé par Kingma & Ba (2014), combine le momentum et le taux d'apprentissage adaptatif. Pour ce dernier, il incorpore une norme infinie pondérée exponentiellement,

$$u_\tau = \max(\omega_2 u_{\tau-1}, \|\mathbf{g}\|). \quad (4.2.20)$$

Son équivalent du momentum est :

$$\mathbf{s}_\tau = \omega_1 \mathbf{s}_{\tau-1} + (1 - \omega_1) \mathbf{g}. \quad (4.2.21)$$

Ainsi, les poids actualisés s'écrivent :

$$\mathbf{w}_\tau = \mathbf{w}_{\tau-1} - \frac{\eta_\tau}{1 - \omega_1} \frac{\mathbf{s}_\tau}{u_\tau}, \quad (4.2.22)$$

où  $\omega_1$  et  $\omega_2$  sont des taux de décroissance.

Une fois que les poids sont mis à jour, le réseau de neurones calcule une nouvelle prédiction avec ses nouveaux paramètres et les étapes se répètent jusqu'à convergence. Puisque la fonction de perte dépend non-linéairement des poids et des biais, elle possède souvent plusieurs minima et autres points stationnaires. Heureusement, pour tirer profit d'un réseau de neurones, le minimum global n'est pas indispensable, d'autant plus qu'on ignore généralement s'il a été atteint (Choromanska et al., 2015; Goodfellow et al., 2014).

#### 4.2.4. Ensembles de données

Afin de s'assurer que le réseau de neurones dénicher les caractéristiques pertinentes, plutôt qu'il apprenne les données par cœur, celles-ci sont réparties entre un ensemble d'entraînement, un ensemble de validation et un ensemble de test. L'apprentissage s'effectue sur l'ensemble d'entraînement qui compte de 60% à 90% des données totales. Simultanément, la performance du réseau de neurones est vérifiée sur l'ensemble de validation. Cette procédure sert à confirmer que la performance se maintient sur des exemples qui n'ont pas guidé l'optimisation. L'ensemble de validation permet aussi de choisir les hyperparamètres, comme le taux d'apprentissage, le nombre de couche cachées, le nombre de neurones par couche, etc. Par ailleurs, pour connaître la performance finale du réseau de neurones entraîné, il faut l'évaluer sur des exemples qui n'ont eu aucune influence sur son apprentissage. L'ensemble de test intervient à cette fin.

### 4.3. Inférence par simulations

La modélisation d'une lentille gravitationnelle est un problème inverse. En d'autres mots, l'image lentillée se détermine facilement à partir d'une source et d'un défecteur connus, mais le problème d'intérêt consiste au contraire à récupérer la source et le défecteur à partir de l'image lentillée, ce qui est plus ardu. La solution réside dans l'inférence par simulations, c'est-à-dire de comparer les observations à des simulations pour déterminer les quantités recherchées, plutôt que de les analyser explicitement. L'intérêt est de s'appuyer sur le problème direct, dont la résolution est simple, tout en profitant des simulateurs complexes et de la puissance de calcul disponibles. Nombre d'inférences menées par divers domaines scientifiques utilisent cette technique.

Toutefois, la comparaison entre les simulations et les observations est un procédé itératif qui se répète jusqu'à ce que l'optimisation converge aux paramètres qui génèrent la simulation la plus fidèle à l'observation. Cette méthode exige beaucoup de temps et de ressources de calcul. D'ailleurs, l'entièreté du processus est reproduit pour chaque observation, ce qui le rend impraticable pour de larges ensembles de données. Aussi, la qualité de l'inférence dépend d'un critère de convergence arbitraire. Surtout, la comparaison entre les simulations et les observations impliquent des assertions sur la distribution de ces dernières et sur celle des paramètres. La raison est que l'inférence s'inscrit dans un cadre statistique bayésien, lequel repose lourdement sur la fonction de vraisemblance, mais que celle-ci n'a pas de forme fermée. (Cranmer et al., 2020)

### 4.3.1. Statistiques bayésiennes

Les statistiques bayésiennes se prêtent bien aux problèmes physiques et scientifiques en général puisqu'elles traitent les paramètres  $\boldsymbol{\pi}$  d'un modèle comme des variables aléatoires, au même titre que les observations  $\mathbf{x}$ . La formule de Bayes consiste à :

$$p(\boldsymbol{\pi} | \mathbf{x}) = \frac{p(\mathbf{x} | \boldsymbol{\pi}) p(\boldsymbol{\pi})}{\int d\boldsymbol{\pi}' p(\mathbf{x} | \boldsymbol{\pi}') p(\boldsymbol{\pi}')}, \quad (4.3.1)$$

où  $p(\boldsymbol{\pi} | \mathbf{x})$  est la distribution a posteriori, soit la probabilité que les paramètres possèdent une valeur spécifique étant donné l'acquisition de ces observations;  $p(\mathbf{x} | \boldsymbol{\pi})$  est la fonction de vraisemblance, soit la probabilité définie par le modèle de bruit qu'une telle observation ait lieu sachant la valeur des paramètres; et  $p(\boldsymbol{\pi})$  est la distribution a priori, c'est-à-dire une supposition préalable sur la distribution des paramètres basée sur des connaissances antérieures. Le dénominateur du côté droit s'écrit aussi  $p(\mathbf{x})$  et se nomme l' « évidence ». Il s'agit de la distribution générale des observations. Inférer les paramètres  $\boldsymbol{\pi}$  se résume donc au calcul de la distribution a posteriori selon cette formule.

Dans le cas d'une inférence par simulations, le simulateur lui-même incarne le modèle statistique. Il prend en entrée les paramètres, exécute à partir de ceux-ci des processus physiques et instrumentaux, puis délivre les résultats finaux. L'inférence de  $H_0$  par les lentilles gravitationnelles fortes, comme bien d'autres, ne vise qu'une partie  $\Theta$  de l'ensemble des paramètres d'entrée  $\boldsymbol{\pi}$ . Ceux qui n'ont pas d'intérêt pour l'inférence en question, soient les paramètres de nuisance  $\zeta$ , doivent être intégrés pour obtenir la fonction de vraisemblance,

$$p(\mathbf{x} | \Theta) = \int d\zeta p(\mathbf{x} | \Theta, \zeta). \quad (4.3.2)$$

La complexité du simulateur et le nombre élevé de dimensions à intégrer rendent cette opération coûteuse en calculs.

### 4.3.2. Estimateur neuronal de ratio

Les estimateurs neuronaux de ratio (NREs) comptent parmi les nouvelles approches, empruntées aux modèles probabilistes en apprentissage automatique, qui permettent à l'inférence par simulations de surmonter l'obstacle de la fonction de vraisemblance complexe. Les NREs sont des réseaux de neurones artificiels qui apprennent à estimer le ratio entre deux distributions. Comme ils servent généralement à des inférences, ces distributions sont souvent la vraisemblance  $p(\mathbf{x}|\Theta)$  et l'évidence  $p(\mathbf{x})$  car le produit de leur ratio avec la distribution a priori donne celle a posteriori. Ce ratio  $r(\mathbf{x}|\Theta)$  s'exprime de deux façons :

$$r(\mathbf{x}|\Theta) \equiv \frac{p(\mathbf{x}|\Theta)}{p(\mathbf{x})} = \frac{p(\mathbf{x}, \Theta)}{p(\mathbf{x})p(\Theta)}. \quad (4.3.3)$$

Comme mentionné précédemment, la vraisemblance  $p(\mathbf{x}|\Theta)$  et l'évidence  $p(\mathbf{x})$  s'obtiennent difficilement. En revanche, il est possible de tirer des paramètres et des observations de la distribution conjointe des paramètres et des observations  $p(\mathbf{x}, \Theta)$  ainsi que du produit de la distribution marginale des paramètres et celle des observations  $p(\mathbf{x})p(\Theta)$ .

L'entraînement des NREs suit donc celui d'une tâche de classification binaire où les deux classes correspondent aux distributions  $p(\mathbf{x}, \Theta)$  et  $p(\mathbf{x})p(\Theta)$ . La classe ayant l'étiquette  $y = 1$  se compose de paires de paramètres  $\Theta$  et de simulations  $\mathbf{x}$  tirées de  $(\mathbf{x}, \Theta) \sim p(\mathbf{x}, \Theta)$ . En d'autres mots, les paires de cette classe sont dépendantes car leurs paramètres sont bel et bien ceux qui ont produit leurs simulations. La classe étiquetée  $y = 0$ , au contraire, contient des paires de paramètres et de simulations indépendantes, puisqu'elles sont formées aléatoirement. Ces paires sont donc tirées de  $(\mathbf{x}, \Theta) \sim p(\mathbf{x})p(\Theta)$ .

La fonction de perte est la même que pour une classification binaire, soit l'entropie croisée. Hermans et al. (2020) obtiennent l'expression du réseau de neurones optimal  $\mathbf{d}^*$  qui minimise la fonction de perte. L'entropie croisée  $\mathcal{L}$  dans le contexte présent s'écrit :

$$\mathcal{L}[\mathbf{d}(\mathbf{x}, \Theta)] = \int d\Theta \int d\mathbf{x} \int d\Theta' p(\Theta) p(\mathbf{x}|\Theta) p(\Theta') [-\log \mathbf{d}(\mathbf{x}, \Theta) - \log(1 - \mathbf{d}(\mathbf{x}, \Theta'))] \quad (4.3.4)$$

$$= \int d\Theta \int d\mathbf{x} p(\Theta) p(\mathbf{x}|\Theta) [-\log \mathbf{d}(\mathbf{x}, \Theta)] + p(\Theta) p(\mathbf{x}) [-\log(1 - \mathbf{d}(\mathbf{x}, \Theta))] . \quad (4.3.5)$$

L'intégrand de l'équation (4.3.5) sera désigné par  $\mathcal{F}(\mathbf{d})$  pour la suite. Le réseau de neurones optimal  $\mathbf{d}^*$  minimise  $\mathcal{F}(\mathbf{d})$ , ce qui signifie que la dérivée de  $\mathcal{F}(\mathbf{d})$  est nulle lorsqu'elle est évaluée à  $\mathbf{d}^*$ .

$$\left. \frac{\partial \mathcal{F}}{\partial \mathbf{d}} \right|_{\mathbf{d}^*} = p(\Theta) p(\mathbf{x}|\Theta) \left[ -\frac{1}{\mathbf{d}^*(\mathbf{x}, \Theta)} \right] + p(\Theta) p(\mathbf{x}) \left[ \frac{1}{1 - \mathbf{d}^*(\mathbf{x}, \Theta)} \right] = 0 \quad (4.3.6)$$

En assumant que  $p(\Theta) > 0$ ,

$$p(\mathbf{x}|\Theta) \frac{1}{\mathbf{d}^*(\mathbf{x}, \Theta)} = p(\mathbf{x}) \frac{1}{1 - \mathbf{d}^*(\mathbf{x}, \Theta)}, \quad (4.3.7)$$

d'où

$$\mathbf{d}^*(\mathbf{x}, \Theta) = \frac{p(\mathbf{x}|\Theta)}{p(\mathbf{x}|\Theta) + p(\mathbf{x})}. \quad (4.3.8)$$

Le ratio défini à l'équation (4.3.3) s'exprime en fonction de  $\mathbf{d}^*(\mathbf{x}, \Theta)$  selon :

$$r(\mathbf{x}|\Theta) = \frac{\mathbf{d}^*(\mathbf{x}, \Theta)}{1 - \mathbf{d}^*(\mathbf{x}, \Theta)}. \quad (4.3.9)$$

Un NRE est donc un estimateur de  $r(\mathbf{x}|\Theta)$ , dénoté par  $\hat{r}(\mathbf{x}|\Theta)$ . Un estimateur de la distribution a posteriori  $\hat{p}(\Theta|\mathbf{x})$  se construit aisément à partir du NRE en le multipliant par une distribution a priori,

$$\hat{p}(\Theta|\mathbf{x}) = p(\Theta) \hat{r}(\mathbf{x}|\Theta). \quad (4.3.10)$$

Les NREs se démarquent par plusieurs avantages notables. D'abord, ils sont amortis, c'est-à-dire que leurs étapes coûteuses en temps et en calculs, comme les simulations et l'entraînement, n'ont pas à être reproduites pour chaque observation. Une fois entraîné, le même NRE peut mener l'inférence sur toutes les observations avec une performance similaire.

Ensuite, un NRE intègre implicitement sur les paramètres de nuisance  $\zeta$  au cours de son entraînement. En effet, la classe d'une paire  $(\mathbf{x}, \Theta)$  est déterminée par les paramètres d'intérêt  $\Theta$ , mais la simulation  $\mathbf{x}$  est tout de même conçue en tirant des paramètres  $\zeta$ . À force de rencontrer des simulations issues de  $\zeta$  dans les deux classes, le NRE apprend à les traiter d'une façon équivalente à une marginalisation.

Finalement, un NRE n'assume ni la forme des distributions impliquées, ni celle de leur ratio. Cette flexibilité lui permet d'estimer avec précision les distributions complexes engendrées par tous les effets inscrits dans le simulateur. Qui plus est, la convergence d'un NRE au véritable ratio des distributions est mathématiquement garantie. Ces deux caractéristiques lui confèrent une grande puissance d'inférence.

## 4.4. Set Transformer

Au même titre que le perceptron multicouche, un Transformer est une architecture de réseau de neurone. Il a été conçu par Vaswani et al. (2017) pour la traduction de textes, c'est-à-dire pour traiter des séquences. Il utilise exclusivement un mécanisme nommé « attention » pour déterminer la corrélation entre deux ou plusieurs éléments de la séquence, comme les mots d'une phrase. Puisqu'un groupe de mots ne se traduit pas nécessairement de la même façon que les mots individuels qui le composent, le Transformer a largement surpassé ses

prédécesseurs.

Lee et al. (2019) ont créé le Set Transformer, soit une variante du Transformer de base qui s'adresse plutôt aux ensembles de nombres. Ces derniers ont une dimension variable et l'ordre de leurs éléments n'a pas d'importance. Un réseau de neurones qui analyse ce type d'entrée devrait donc supporter des entrées de différentes tailles. De plus, sa sortie devrait être invariante sous la permutation des éléments d'un ensemble. Le Set Transformer utilise le mécanisme d'attention pour assembler des blocs d'opération qui sont équivariants ou invariants sous la permutation des éléments d'un ensemble. Lee et al. (2019) ont aussi prouvé que le Set Transformer est un approximateur universel des fonctions invariantes sous permutation.

#### 4.4.1. Attention

Considérons deux ensembles, soient  $X$  et  $Y$ , contenant chacun  $m$  éléments de dimension  $d$ . Pour évaluer l'attention entre eux, on assigne  $X$  aux *queries*  $\mathcal{Q} \in \mathbb{R}^{m \times d}$  et  $Y$  aux *keys*  $\mathcal{K} \in \mathbb{R}^{m \times d}$  ainsi qu'aux *values*  $\mathcal{V} \in \mathbb{R}^{m \times d}$ . L'attention s'écrit alors :

$$\text{Attention}(\mathcal{Q}, \mathcal{K}, \mathcal{V}) = \sigma\left(\frac{\mathcal{Q}\mathcal{K}^\top}{\sqrt{d}}\right) \mathcal{V}, \quad (4.4.1)$$

où  $\sigma(\cdot)$  est la fonction softmax,

$$\sigma_i(\mathbf{x}) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}. \quad (4.4.2)$$

Le produit  $\mathcal{Q}\mathcal{K}^\top \in \mathbb{R}^{m \times m}$  quantifie la similarité entre les paires de *queries* et de *keys*. La fonction softmax contraint les éléments entre 0 et 1, et ce, de façon à ce que leur somme soit unitaire. Cette opération équivaut à calculer les probabilités de classes mutuellement exclusives. Le résultat final correspond à une somme pondérée des *values*  $\mathcal{V}$  où le poids est d'autant plus grand que le produit scalaire entre une *key* et une *query* est élevé.

#### 4.4.2. Attention à têtes multiples

Au lieu d'évaluer l'attention une seule fois, Vaswani et al. (2017) projettent  $\mathcal{Q}$ ,  $\mathcal{K}$  et  $\mathcal{V}$  sur  $h$  têtes de dimension  $d/h$ , calculent l'attention pour chacune d'elles, concatènent les résultats, puis les transforment linéairement. Mathématiquement, cela se traduit par :

$$\text{Multihead}(\mathcal{Q}, \mathcal{K}, \mathcal{V}) = \text{concatenate}(O_1, \dots, O_h) \mathbf{w}^O, \quad (4.4.3)$$

où :

$$O_j = \text{Attention}\left(\mathcal{Q}\mathbf{w}_j^{\mathcal{Q}}, \mathcal{K}\mathbf{w}_j^{\mathcal{K}}, \mathcal{V}\mathbf{w}_j^{\mathcal{V}}\right). \quad (4.4.4)$$

Les poids de ce réseau de neurones comprennent donc  $\{\mathbf{w}_j^{\mathcal{Q}}, \mathbf{w}_j^{\mathcal{K}}, \mathbf{w}_j^{\mathcal{V}}\}_{j=1}^h$ , où  $\mathbf{w}_j^{\mathcal{Q}}, \mathbf{w}_j^{\mathcal{K}}, \mathbf{w}_j^{\mathcal{V}} \in \mathbb{R}^{d \times d/h}$ , et  $\mathbf{w}^O \in \mathbb{R}^{d \times d}$ .

Lee et al. (2019) définissent un bloc d'attention à têtes multiples de la façon suivante :

$$\text{MAB}(X, Y) = \text{LayerNorm}(L + \text{rFF}(L)) , \quad (4.4.5)$$

avec :

$$L = \text{LayerNorm}(X + \text{Multihead}(X, Y, Y)) , \quad (4.4.6)$$

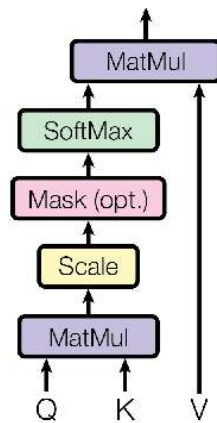
où  $\text{LayerNorm}$  est une normalisation par couche (Ba et al., 2016) et  $\text{rFF}$  est une couche de neurones.

L'opérateur SAB calcule plutôt l'auto-attention à têtes multiples, soit l'attention d'un ensemble  $X$  avec lui même. Il se résume à :

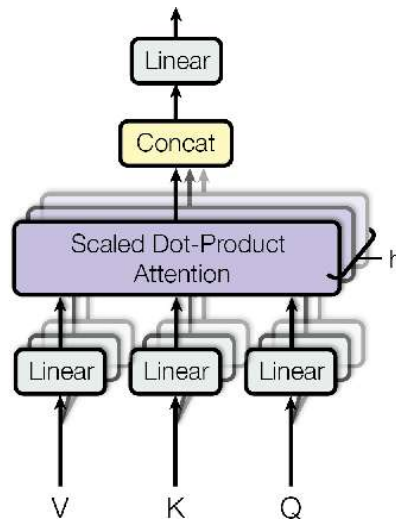
$$\text{SAB}(X) = \text{MAB}(X, X) . \quad (4.4.7)$$

Comme souhaité, cet opérateur est équivariant sous permutation des éléments de  $X$ .

led Dot-Product Attention



Multi-Head Attention



**Fig. 4.5.** À gauche, schéma de la fonction d'attention. À droite, schéma de l'attention à têtes multiples. Crédit : Vaswani et al. (2017)

### 4.4.3. Architecture

Le Set Transformer se compose d'un encodeur et d'un décodeur. L'encodeur prend en entrée un ensemble  $X$  et lui applique successivement des blocs d'auto-attention.

$$\text{Encoder}(X) = \text{SAB}(\text{SAB}(X)) \quad (4.4.8)$$

La première étape du décodeur consiste à réduire l'information extraite de chaque ensemble à une seule quantité. Pour ce faire, le décodeur passe la sortie  $Z$  de l'encodeur à des couches de neurones, puis calcule l'attention à têtes multiples entre le résultat et un vecteur de poids



$S$ .

$$\text{PMA}(Z) = \text{MAB}(S, \text{rFF}(Z)) \quad (4.4.9)$$

Cette transformation remplit les critères nécessaires au traitement d'ensembles, puisqu'elle est invariante sous permutation. De plus, la dimension de ses entrées n'affecte pas celle de son résultat. Par la suite, le décodeur applique à nouveau l'auto-attention. Il se termine par des couches de neurones qui fournissent la prédiction.

$$\text{Decoder}(Z) = \text{rFF}(\text{SAB}(\text{PMA}(Z))) \quad (4.4.10)$$

Cette architecture convient bien au traitement de sources ponctuelles fortement lentillées. En effet, leurs images produites par un système de lentille gravitationnelle forment un ensemble, c'est-à-dire que l'ordre dans lequel elles sont considérées n'a pas d'impact sur l'inférence. Le Set Transformer intègre explicitement cette propriété à son analyse. De plus, un même Set Transformer peut traiter à la fois les ensembles de deux images et ceux de quatre images. Cette architecture est donc toute désignée pour un NRE entraîné à inférer la constante de Hubble à partir de lentilles gravitationnelles fortes.



## Chapitre 5

---

# Time Delay Cosmography with a Neural Ratio Estimator

Ève Campeau-Poirier,<sup>1,2</sup> Laurence Perreault Levasseur,<sup>1,2,3</sup> Yashar Hevazeh<sup>1,3</sup> Adam Coogan<sup>1,2</sup>

<sup>1</sup>*Département de physique, Université de Montréal, Montréal, H3C 3J7, Canada*

<sup>2</sup>*Mila - Quebec Artificial Intelligence Institute, Montréal, Canada*

<sup>3</sup>*Center for Computational Astrophysics, Flatiron Institute, 162 5th Avenue, 10010, New York, NY, USA*

Cet article sera soumis à la revue *The Astrophysical Journal* (ApJ) dans les prochains mois.

## Résumé

Nous explorons l'utilisation d'un estimateur neuronal de ratio (NRE) pour déterminer la constante de Hubble ( $H_0$ ) à partir de lentilles gravitationnelles fortes. En supposant un profil de masse isotherme ellipsoïde singulier (SIE) pour le défecteur, nous simulons des mesures de délais, des mesures de positions des images et des paramètres de lentille modélisés. Seul du bruit gaussien imite les erreurs et les incertitudes pour garder la vraisemblance traitable, et ainsi, comparer le NRE avec une méthode conventionnelle, soit le *nested sampling*, dans un cas où cette dernière est fiable. Nous entraînon le NRE à évaluer la distribution a posteriori de  $H_0$  en fonction des mesures de délais, des potentiels de Fermat relatifs (calculés à partir des paramètres modélisés et des positions d'images mesurées), du décalage vers le rouge du défecteur et celui de la source. Les résultats du NRE sont près de ceux de la méthode conventionnelle, ils présentent des incertitudes quelque peu surestimées et ils se combinent sans biais lors d'une inférence sur une population.

**Mots-clés:** Constante de Hubble — Lentilles gravitationnelles fortes — Estimateur neuronal de ratio — Inférence par simulations — Apprentissage automatique.

## Abstract

We explore the use of a Neural Ratio Estimator (NRE) to determine the Hubble constant ( $H_0$ ) in the context of time delay cosmography. Assuming a Singular Isothermal Ellipsoid mass profile for the deflector, we simulate time delay measurements, image position measurements, and modeled lensing parameters. Only Gaussian noise emulates the errors and uncertainties to keep the likelihood tractable, and thus compare the NRE with a conventional method, nested sampling, in a case where the latter is reliable. We train the NRE to output the posterior distribution of  $H_0$  given the time delay measurements, the relative Fermat potentials (calculated from the modeled parameters and the measured image positions), the deflector redshift, and the source redshift. The NRE posteriors are close to the ones from the conventional method, they show slightly overestimated uncertainties, and they combine in a population inference without bias.

**Keywords:** Hubble constant — Strong gravitational lensing — Neural ratio estimator — Simulation-based inference — Machine learning.

## 5.1. Introduction

Even after years of careful searches for systematic effects, a  $4\sigma$  to  $6\sigma$  tension persists between the Hubble constant ( $H_0$ ) measurements from early- and late-universe probes (Di Valentino et al., 2021). Time delay cosmography can provide an independent measurement of  $H_0$ , under a different set of assumptions and with different systematics from methods relying on the local distance ladder or on the cosmic microwave background (CMB). This could provide a path for resolving the current crisis in cosmology, by helping us determine if this discrepancy can be attributed to a yet unknown systematics or, alternatively, new physics outside the  $\Lambda$ CDM model.

Time delay cosmography can be used to measure  $H_0$  by determining the time delays between the multiple images of a variable background light source strongly lensed by a foreground mass deflector. The previous determinations of  $H_0$  with this method, conducted on a sample of seven lensing systems, could not achieve a sufficient precision to solve the discrepancy (Birrer, Simon & Treu, Tommaso, 2021). In its ten years of activity, the Legacy Survey of Space and Time (LSST) at the Vera Rubin Observatory is expected to solve this issue by detecting more than 8,000 strongly lensed quasar, among which 3,000 will have well-measured time delays (Oguri & Marshall, 2010), and 500 strongly lensed supernovae (Goldstein & Nugent, 2016).

However, the current analysis to infer  $H_0$  through time delay cosmography involves a time-consuming forward-modeling of the lensed images, which must be repeated for each observation. This procedure is not feasible on a dataset of the size expected from LSST due to computational cost. Therefore, many efforts have been made to accelerate the analysis through machine learning (Hezaveh et al., 2017; Levasseur et al., 2017; Morningstar et al., 2019; Pearson et al., 2019; Schuldt, S. et al., 2021; Park et al., 2021).

Once the lensing parameters have been modeled by an artificial neural network, the likelihood to infer  $H_0$  is still intractable. One must then resort to a Markov-Chain Monte Carlo (MCMC) to evaluate the  $H_0$  posterior, which is a time-consuming method performed on a single lensing system at a time. Another option is to make simplifying assumptions to employ an explicit-likelihood method, but this can lead to a biased inference.

Simulation-based inference frameworks are particularly well-suited for time delay cosmography because they only require an accurate simulation pipeline. Among them, Neural Ratio Estimators (NREs) allow implicit marginalization over large sets of nuisance parameters, while providing an efficient way to estimate low-dimensional variables. Unlike a Bayesian Neural Networks (BNN), NREs do not assume the form of any distribution, which makes them highly flexible. Moreover, their calculations are as fast as any other machine learning algorithm, and the same NRE can perform the inference on different examples with the same

level of accuracy.

In this work, we explore the application of a NRE to time delay cosmography. Because our aim is a proof of concept, we prioritize simplicity over realism. Therefore, we assume a simple model for the deflector’s mass density profile, i.e. a Singular Isothermal Ellipsoid, to simulate time delays and image positions. For the same reason, we use Gaussian noise to mimic the uncertainties yielded by lens modeling and time delay measurements. From the modeled parameters, we compute the Fermat potential at the measured image positions. The NRE learns to predict the  $H_0$  posterior distribution given the calculated Fermat potentials, the time delay measurements, the deflector redshift and the source redshift.

The likelihood being tractable in this case, we compare the NRE predictions with posteriors obtained from nested sampling on different lensing systems and different  $H_0$  values. We also assess the NRE statistical consistency with a coverage diagnostic on 10,000 examples. Finally, we conduct a population inference on sets of 50, 500, 3,000, and 8,000 synthetic lensing systems to explore the NRE’s performance on set sizes of different orders of magnitude, and on those expected by LSST. We find that the NRE agrees with the nested sampling, that it weakly overestimates the uncertainties, and that it leads to unbiased predictions on population inferences.

Section 5.2 explains the method to measure  $H_0$  with time delay cosmography, and the basic principles of a NRE. In section 5.3, we describe the assumptions and procedures to simulate our data. Section 5.4 presents the artificial neural network architecture and training process. In section 5.5, we show the resulting posterior distributions, we test their statistical coverage, and we conduct a population inference. We discuss how to extend this framework and conclude in section 5.6.

## 5.2. Method

Section 5.2.1 provides a quick overview of the time delay cosmography theory and of the current analysis methods. In section 5.2.2, we explain the theory underlying Neural Ratio Estimators, and how we can apply them to time delay cosmography.

### 5.2.1. Time delay cosmography

Gravitational lensing occurs when light rays follow a space-time geodesic curved by matter. In this situation, an observer detects the light rays at positions different from their initial ones, which results in a distorted image. The lensing equation,

$$\boldsymbol{\beta} = \boldsymbol{\theta} - \boldsymbol{\alpha}(\boldsymbol{\theta}) , \tag{5.2.1}$$

summarizes this phenomenon by retracing the source plane angular position  $\boldsymbol{\beta}$  of a ray observed at the image plane angular position  $\boldsymbol{\theta}$  after a mass deflector has deviated it by an angle  $\boldsymbol{\alpha}$ . The lensing potential  $\psi$  of the massive object determines the angular deflection  $\boldsymbol{\alpha}$  according to :

$$\boldsymbol{\alpha}(\boldsymbol{\theta}) = \nabla\psi(\boldsymbol{\theta}) . \quad (5.2.2)$$

It is also related to the dimensionless surface mass density  $\kappa$  as follows :

$$\nabla^2\psi(\boldsymbol{\theta}) = 2\kappa(\boldsymbol{\theta}) . \quad (5.2.3)$$

Furthermore, gravitational lensing affects the light rays travel time from their source to the observer. The presence of a mass deflector in the light's trajectory lengthens its travel time by :

$$t(\boldsymbol{\theta}, \boldsymbol{\beta}) = \frac{D_{\Delta t}}{c} \phi(\boldsymbol{\theta}, \boldsymbol{\beta}) , \quad (5.2.4)$$

where  $D_{\Delta t}$  is the time delay distance,  $c$  is the speed of light, and  $\phi$  is the Fermat potential. The latter is only determined by the lensing system, and is defined as :

$$\phi(\boldsymbol{\theta}, \boldsymbol{\beta}) \equiv \frac{(\boldsymbol{\theta} - \boldsymbol{\beta})^2}{2} - \psi(\boldsymbol{\theta}) . \quad (5.2.5)$$

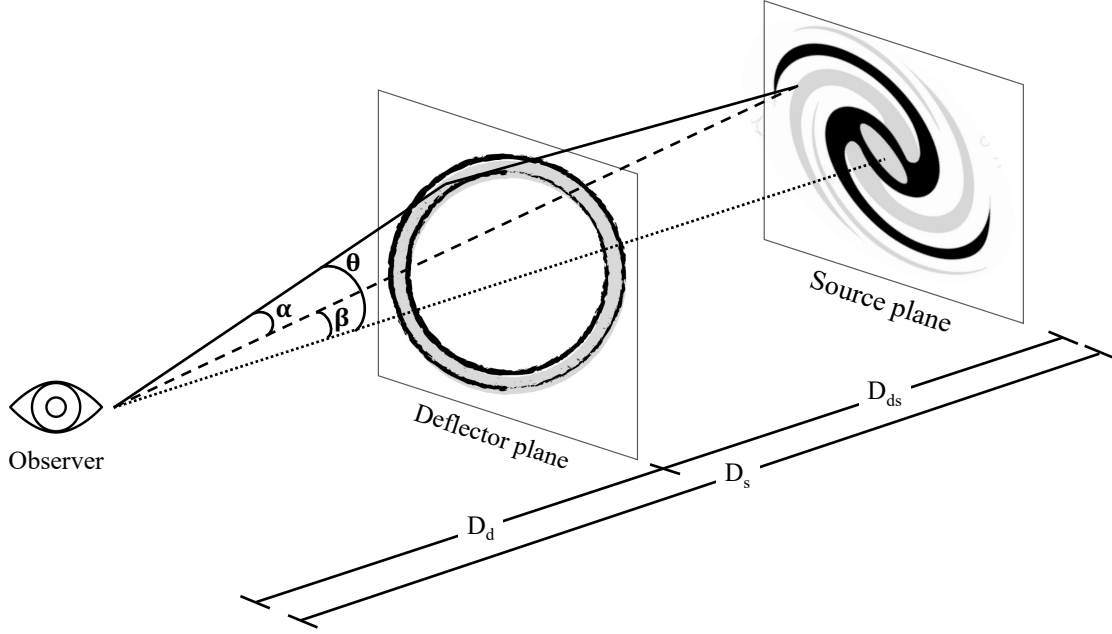
The time delay distance is inversely proportional to  $H_0$ , and weakly depends on other cosmological parameters. It is written :

$$D_{\Delta t} \equiv (1 + z_d) \frac{D_d D_s}{D_{ds}} , \quad (5.2.6)$$

where  $z_d$  is the deflector redshift,  $D_d$  is the diameter angular distance between the observer and the deflector,  $D_s$  is the diameter angular distance between the observer and the source,  $D_{ds}$  is the diameter angular distance between the deflector and the source, see Figure (5.1).

To infer  $H_0$  with time delay cosmography, one must measure the excess of light travel time  $t(\boldsymbol{\theta}, \boldsymbol{\beta})$  and model the Fermat potential  $\phi(\boldsymbol{\theta}, \boldsymbol{\beta})$  in order to isolate the  $H_0$  dependency in  $D_{\Delta t}$ . This can be done by observing a strongly lensed variable light source. Strong gravitational lensing occurs when a light ray's initial position is inside the region where the foreground mass will deflect it on multiple paths, producing as many images of the background source. In most cases, strong gravitational lensing systems yield two or four images, referred to as doubles and quads, respectively.

Each path is affected by a different Fermat potential, resulting in a different light travel time. Therefore, the source fluctuations are distinguishable at different moments on the light curve of each image. This allows to evaluate the relative travel times between paths  $\Delta t$ , which are called time delays. They are calculated between an arbitrarily chosen image and the others.



**Fig. 5.1.** Diagram of gravitational lensing. The solid line illustrates the path taken by the deflected light. The dashed line indicates the path taken by the light in absence of the deflector. The dotted line shows the center of the coordinate system. The difference of travel time between the paths represented by the solid and the dashed lines is expressed by Equation (5.2.4).

The equation to solve for  $H_0$  becomes :

$$\Delta t \equiv \frac{D_{\Delta t}}{c} \Delta \phi, \quad (5.2.7)$$

where

$$\Delta \phi = \phi(\boldsymbol{\theta}_A, \boldsymbol{\beta}) - \phi(\boldsymbol{\theta}_B, \boldsymbol{\beta}), \quad (5.2.8)$$

with A and B designating the image of reference and another one.

In this framework, the posterior distribution of  $H_0$  takes generally the form :

$$P(H_0 | \Delta \mathbf{t}, \mathbf{d}) \propto \int d\boldsymbol{\zeta} P(\Delta \mathbf{t} | H_0, \boldsymbol{\zeta}, \mathbf{M}) P(\boldsymbol{\zeta} | \mathbf{d}, \mathbf{M}) P(H_0), \quad (5.2.9)$$

where  $\mathbf{d}$  is the lensing observation and  $\boldsymbol{\zeta}$  is a set of parameters describing the lensing system.  $\mathbf{M}$  includes all observational effects (e.g. instrumental noise, point spread function, image covariance matrix, deflector's light, dust, etc.), which make the likelihood intractable. This problem is often circumvented with approximations, such as assuming a normal distribution. However, this work-around may introduce biases. In addition, the integral is computed with a MCMC algorithm, which is a tedious and time-consuming procedure.



### 5.2.2. Neural Ratio Estimator

A forward model that generates synthetic observations can incorporate the observational effects  $\mathbf{M}$  mentioned in section 5.2.1. This makes simulation-based inference well-suited for time delay cosmography. In this approach, a simulator becomes the statistical model : samples can be drawn from it, even though the probability densities are intractable. The posterior distribution can thus be learned from the samples, with an accuracy mostly limited by the realism of the simulations.

A Neural Ratio Estimator (Cranmer et al., 2015) is one of the strategy to learn the posterior distribution of parameters  $\Theta$  given samples  $\mathbf{x}$  from simulations. It is a neural network that estimates the ratio between two distributions. Its core is a discriminator trained to distinguish samples drawn from both of them. A NRE is therefore convenient to estimate a posterior distribution, because there are two ways to express it as a ratio :

$$p(\Theta | \mathbf{x}) = \frac{p(\Theta)p(\mathbf{x} | \Theta)}{p(\mathbf{x})} = \frac{p(\Theta)p(\mathbf{x}, \Theta)}{p(\mathbf{x})p(\Theta)}. \quad (5.2.10)$$

With the simulator, we draw dependent sample-parameter pairs from the joint distribution  $p(\mathbf{x}, \Theta)$ . The class label  $y = 1$  is assigned to them. The second distribution is the product of the sample and the parameter marginal distributions  $p(\mathbf{x})p(\Theta)$ . The sample-parameter pairs drawn from this distribution are independent, and have the class label  $y = 0$ .

The optimal discriminator  $\mathbf{d}^*$  that classifies samples from these two distributions converges to the decision function,

$$\mathbf{d}^*(\mathbf{x}, \Theta) = p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x}, \Theta)}{p(\mathbf{x}, \Theta) + p(\mathbf{x})p(\Theta)}. \quad (5.2.11)$$

The ratio  $r(\mathbf{x} | \Theta)$  between the distributions can be written as a function of the discriminator :

$$r(\mathbf{x} | \Theta) \equiv \frac{p(\mathbf{x}, \Theta)}{p(\mathbf{x})p(\Theta)} = \frac{\mathbf{d}^*(\mathbf{x}, \Theta)}{1 - \mathbf{d}^*(\mathbf{x}, \Theta)}. \quad (5.2.12)$$

The product between the estimator  $\hat{r}(\mathbf{x} | \Theta)$  and the prior distribution gives a posterior distribution estimator because,

$$\begin{aligned} p(\Theta)r(\mathbf{x} | \Theta) &= p(\Theta) \frac{p(\mathbf{x}, \Theta)}{p(\mathbf{x})p(\Theta)} \\ &= p(\Theta) \frac{p(\mathbf{x} | \Theta)}{p(\mathbf{x})} \\ &= p(\Theta | \mathbf{x}). \end{aligned} \quad (5.2.13)$$

To conduct an inference with a trained Neural Ratio Estimator, the ratio  $r(\mathbf{x} | \Theta)$  is calculated multiple times for the same observation, but with a different parameter values at each computation. Multiplied by the prior distribution, this gives an estimate of the posterior

distribution on a grid of parameter values.

For time delay cosmography, the parameters  $\Theta$  to infer boil down to  $H_0$ , but there are many choices of observables  $\mathbf{x}$ . They could be the raw data (image light curves and the lensing observation), analysis products (time delay measurements, inferred lensing parameters or modeled Fermat potentials), or a combination of both. Because our work is a proof of principle, we choose the option that provides the most explicit information to the NRE, and that has the smallest dimensionality, which is to input the time delay measurements and their corresponding Fermat potential estimates. This choice is also compatible with both lens modeling with a parametric model and free-form reconstruction.

This approach improves upon the traditional analysis in many ways. The NRE is highly flexible because it learns the ratio of distributions without any assumption about their form. This resolves the issue of the intractable likelihood without using approximations. Also, the NRE replaces the MCMC for the evaluation of the  $H_0$  posterior distributions, which is more efficient because the NRE is amortized. Furthermore, not only is the NRE inference on a single system fast, it can be done simultaneously on a large number due to its leverage of GPUs computing power.

### 5.3. Simulations

In our experiment, the lensing systems consist of a variable point source, and a deflector with a Singular Isothermal Ellipsoid (SIE ; Kormann et al., 1994) mass profile. This parametric model assumes that a galaxy behaves as an ideal gas in hydrostatic and thermal equilibrium, where the stars act like particules, and confined by an axisymmetric self-gravitational potential. The dimensionless surface mass density of a SIE is :

$$\kappa(\boldsymbol{\theta}) = \frac{\theta_E}{2} \sqrt{\frac{f}{f^2\theta_1^2 + \theta_2^2}}, \quad (5.3.1)$$

where  $\theta_1$  and  $\theta_2$  are the two components of the angular position  $\boldsymbol{\theta}$ ,  $f$  is the axis ratio, and  $\theta_E$  is the Einstein radius. Despite that this expression diverges at the origin and that it implies an infinite mass, it is a good approximation of the lens potential on galaxy scales.

According to Equation (5.2.5), computing the Fermat potential at the image positions requires the lensing potential, the source position, and the image positions. The lensing potential for the SIE is obtained from Equations (5.2.3) and (5.3.1). We also add external shear with modulus  $\gamma_{\text{ext}}$  and orientation  $\varphi_{\text{ext}}$ .

We fix the source position  $(x_s, y_s)$  at the center of the field of view. As for the image positions, the SIE profile offers a simple procedure to solve the lensing equation for them. One

Parameter	Distribution
<b>Cosmology</b>	
Hubble constant (km s <sup>-1</sup> Mpc <sup>-1</sup> )	$H_0 \sim \text{U}(50, 90)$
Dark energy density	$\Omega_\Lambda = 0.7$
Matter energy density	$\Omega_m = 0.3$
<b>Deflector</b>	
Redshift	$z_d \sim \text{U}(0.04, 0.5)$
Position (")	$x_d, y_d \sim \text{U}(-0.3, 0.3)$
Einstein radius (")	$\theta_E \sim \text{U}(0.5, 2.0)$
Axis ratio	$f \sim \text{U}(0.30, 0.99)$
Orientation (rad)	$\varphi_d \sim \text{U}(-\pi/2, \pi/2)$
<b>External Shear</b>	
Modulus	$\gamma_{\text{ext}} \sim \text{U}(0, 0.2)$
Orientation (rad)	$\varphi_{\text{ext}} \sim \text{U}(-\pi/2, \pi/2)$
<b>Variable point light source</b>	
Redshift	$z_s \sim \text{U}(1, 3)$
Position (")	$x_s, y_s = (0, 0)$

**Table 5.1.** Prior distributions of all the parameters needed to generate Fermat potentials and time delays in our framework

has to find the roots of :

$$\left[ \beta_1 + \frac{\sqrt{f}}{\sqrt{1-f^2}} \arcsin\left(\frac{\sqrt{1-f^2}}{f} \cos \vartheta\right) \right] \sin \vartheta - \left[ \beta_2 + \frac{\sqrt{f}}{\sqrt{1-f^2}} \arcsin\left(\sqrt{1-f^2} \sin \vartheta\right) \right] \cos \vartheta, \quad (5.3.2)$$

where  $\vartheta$  is the orientation of the vector  $\boldsymbol{\theta}$  in the deflector's plane. The conversion of the results into the field of view depend on the SIE central position  $(x_d, y_d)$  and on its orientation  $\varphi_d$ . Thus, we sample seven parameters from uniform distributions to generate the Fermat potentials :  $x_d, y_d, \theta_E, f, \varphi_d, \gamma_{\text{ext}}$ , and  $\varphi_{\text{ext}}$ . Table 5.1 list all the prior distributions.

We compute time delay distances according to Equation (5.2.6). The  $H_0$  value, the source redshift and the deflector redshift are drawn from uniform prior distributions detailed in Table (5.1). We assume a flat  $\Lambda$ CDM cosmological model with matter and dark energy densities  $\Omega_m = 0.3$  and  $\Omega_\Lambda = 0.7$ . With the Fermat potential at the image positions and the time delay distance, we calculate the excess travel time from Equation (5.2.4). We then select the minimum value as the reference one, and subtract it from the others to obtain the time delays of Equation (5.2.7). This means doubles have one time delay - Fermat potential pair, and quads have three of them, all positive.

To train the NRE, we want to emulate the results of a standard analysis, which models the system parameters from the lensing observation and measures the time delays from the image light curves. Therefore, we add Gaussian noise to the lensing parameters and image positions to mimic the analysis products. We again compute the Fermat potentials, but with

<b>Noise standard deviation</b>	
<b>Observables</b>	
Time delays	0.35
Image positions	$5 \times 10^{-4}$
<b>Deflector</b>	
Position (")	$5 \times 10^{-4}$
Einstein radius (")	$5 \times 10^{-4}$
Axis ratio	$5 \times 10^{-4}$
Orientation (rad)	$5 \times 10^{-4}$
<b>External shear</b>	
Modulus	$5 \times 10^{-4}$
Orientation (rad)	$5 \times 10^{-4}$
<b>Active galactic nucleus</b>	
Position (")	$5 \times 10^{-4}$

**Table 5.2.** Standard deviation of the Gaussian noise distributions used to mimic the uncertainties of lens modeling, time delay measurements, and image position measurements

these noisy parameters. These Fermat potentials are the ones that the NRE takes as inputs. For the time delays, we add Gaussian noise to the ones generated with the true parameters before feeding them to the NRE. This replicates the uncertainty yielded by the light curve measurements. Table (5.2) summarizes all the standard deviations of the Gaussian noise distributions. They were chosen to yield uncertainties similar to those of lens modeling.

We also give the redshifts to the NRE because they influence the time delay distance, but the Fermat potentials do not contain any information about them. Thus, they are necessary to discriminate the dependent and independent  $H_0$ - $\Delta t$  pairs. Moreover, because they are easily measured, we consider that they are exactly known. This framework makes the true posteriors tractable, allowing the assessment of the NRE’s accuracy.

## 5.4. Neural Network

The network architecture is presented in section 5.4.1. Section 5.4.2 provides the training procedure for the NRE.

### 5.4.1. Architecture

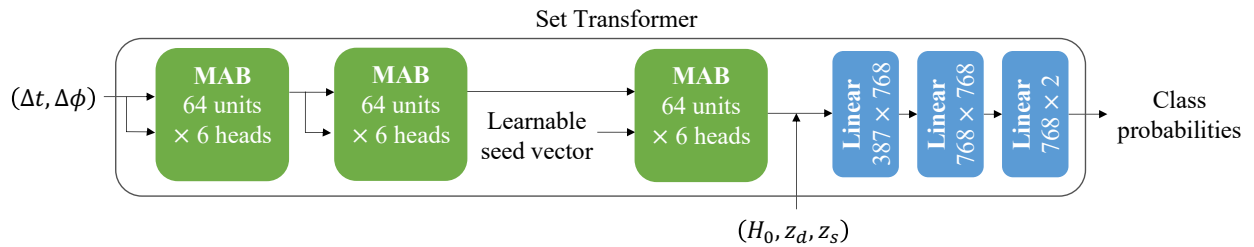
The data structure motivates our choice of architecture. Each image is characterized by a pair of time delay - relative Fermat potential. Each double and each quad is represented by one pair and three pairs, respectively. The images of a lensing system form a set, because their order has no meaning, i.e. it can be chosen arbitrarily and it should not affect the inference result. In addition, the presence of doubles and quads in the dataset implies that the inputs have different sizes. Therefore, the neural network must be invariant under the

permutation of set elements, and it must accept inputs of variable sizes.

The Deep Sets model (Zaheer et al., 2017) meets these criteria by computing features for all set elements, then performing a permutation invariant pooling operation, e.g. a sum or a mean, across set elements. Another architecture, the Set Transformer (Lee et al., 2019), meets these requirements by applying self-attention to the set. We experimented with both, but we report only the results obtained with the Set Transformer because it outperformed the Deep Sets model.

Figure 5.2 illustrates our Set Transformer architecture. The first self-attention block computes multi-head attention between the time delay - relative Fermat potential pairs belonging to the same lensing system. The second self-attention block repeats the operation with the output of the first one. After, the features are aggregated by computing multi-head attention between a learnable seed vector and them. At each step, we use 6 attention heads of dimension 64. The  $H_0$  value,  $z_d$  and  $z_s$  are concatenated to the result, which is then fed sequentially to 3 linear layers, each of 768 neurons. There is a ELU activation functions before and after the second layer. The whole neural network counts 2,224,514 parameters.

At inference time, we apply a softmax function the final output to retrieve the class probabilities. We then insert the probability of the class with label  $y = 1$  in Equation (5.2.12) to estimate the distribution ratio. The latter is equivalent to the posterior density at the input  $H_0$  because the prior is uniform.



**Fig. 5.2.** Set Transformer architecture for the discriminator. The green squares represent the multihead attention blocks (MAB), and the blue rectangles, the linear layers. There are ELU activation functions after the first and the second linear layers. At inference time, a softmax function is added after the last layer. See Table A.1 in Appendix A for the input sizes at each step.

### 5.4.2. Training

We generate 400,000 examples, then split them between the training set, the validation set and the test set with conventional proportions, i.e. 80%, 10%, and 10%, respectively. We remove the examples with  $H_0$  outside the interval between  $65 \text{ km Mpc}^{-1} \text{ s}^{-1}$  and 75

km Mpc<sup>-1</sup> s<sup>-1</sup> from the test set to ensure that the inference is not affected by border effects. Therefore, the training set, the validation set and the test set contain respectively 320,000 examples, 40,000 examples and 10,000 examples. All the datasets are composed of approximately 45% doubles and 55% quads. An example is composed of the time delays of a lensing system, its parameters, the source and deflector redshifts, and the  $H_0$  value that produced the time delays. We train the neural network on batches of 256 examples with a binary cross entropy loss as the objective function. At each batch, we draw a new realization of noise for the time delays and for the parameters. We then compute the Fermat potentials. After, we select half of the examples, we exchange their  $H_0$  values, and we assign class label 0 to them. The training lasts for 5,000 epochs. The learning rate starts at  $1 \times 10^{-4}$ , and decreases by half every 500 epochs. This schedule is the optimal one found with a hyper parameter search.

## 5.5. Results

The posterior of our experiment can be written as :

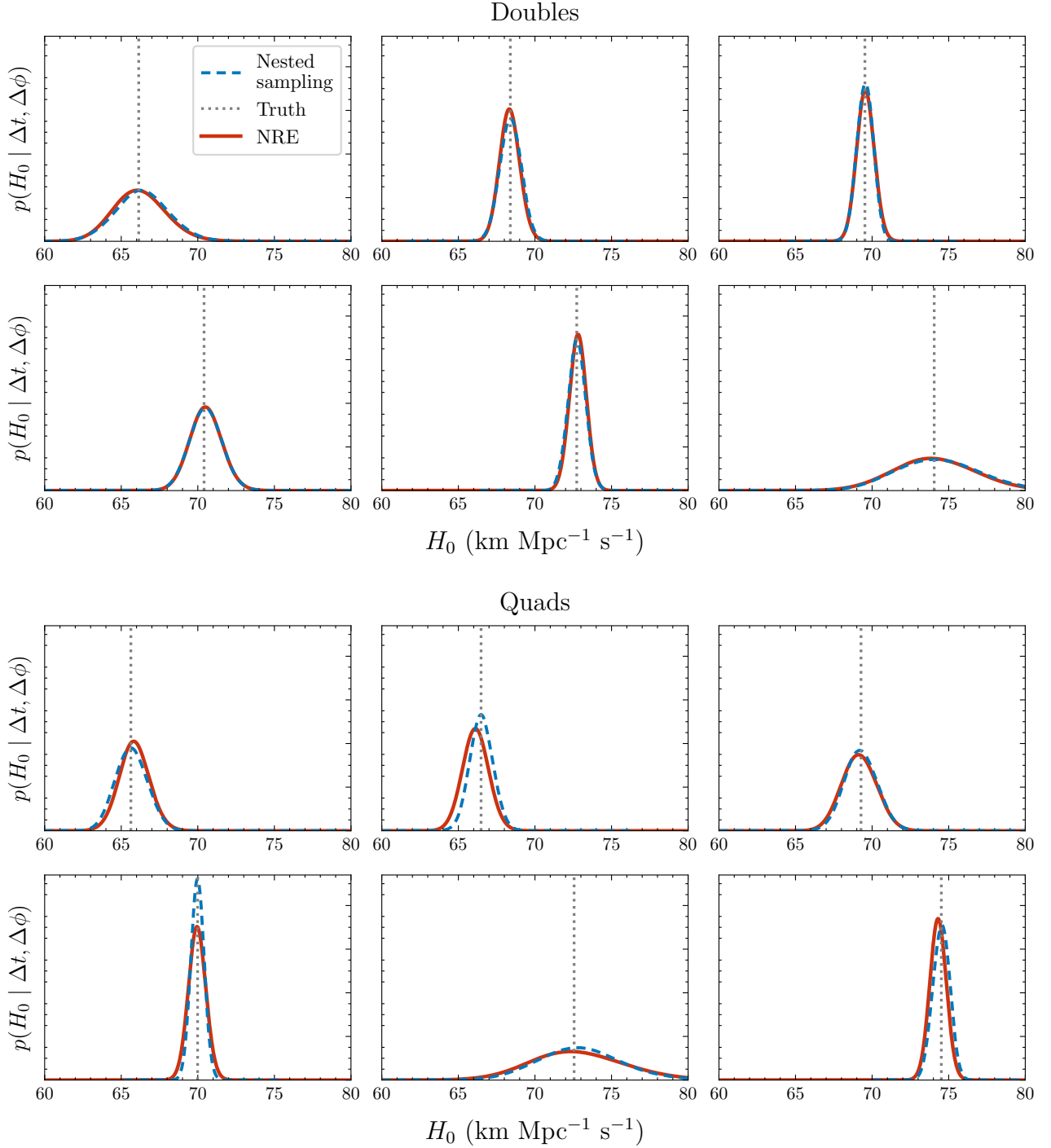
$$P(H_0|\Delta t, \Delta\phi) \propto P(\Delta t|H_0, \Delta\phi)P(\Delta\phi)P(H_0), \quad (5.5.1)$$

with

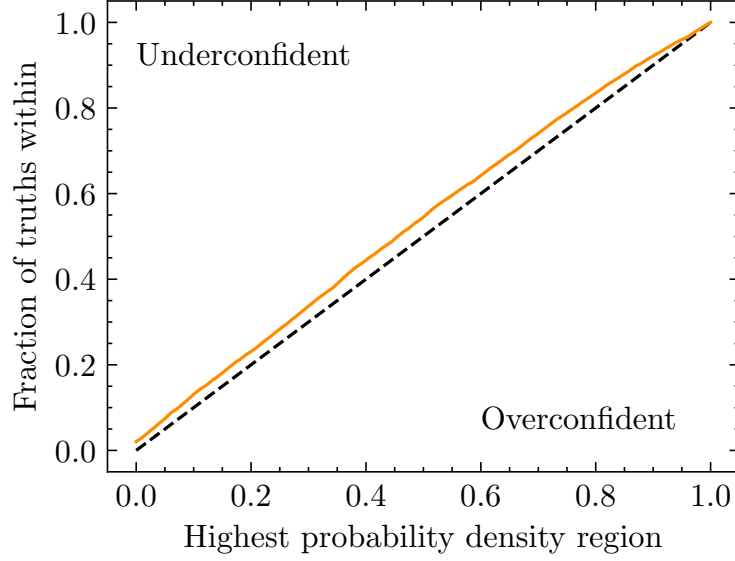
$$P(\Delta\phi) \propto \int d\zeta P(\Delta\phi|\zeta)P(\zeta), \quad (5.5.2)$$

where  $P(\Delta t|H_0, \Delta\phi)$  and  $P(\zeta)$  are normal distributions,  $P(\Delta\phi|\zeta)$  is a delta function, and  $P(H_0)$  is a uniform distribution. In Figure 5.3, we compare the NRE results on 12 test examples with those of nested sampling performed with the package POLYCHORD (Handley et al., 2015a,b). Each plot is associated to a different lensing system and a different  $H_0$  value. The nested sampling and the NRE posteriors are respectively indicated by the blue dashed line and the red solid line. The NRE shows little discrepancy with the nested sampling posteriors on doubles. Even if they remain reasonable, the dissimilarities between the two distributions are more noticeable on quads. For most quads, when there is an offset between the NRE distribution and the truth, the NRE distribution is also wider than the nested sampling one, which mean it is less precise, but still accurate.

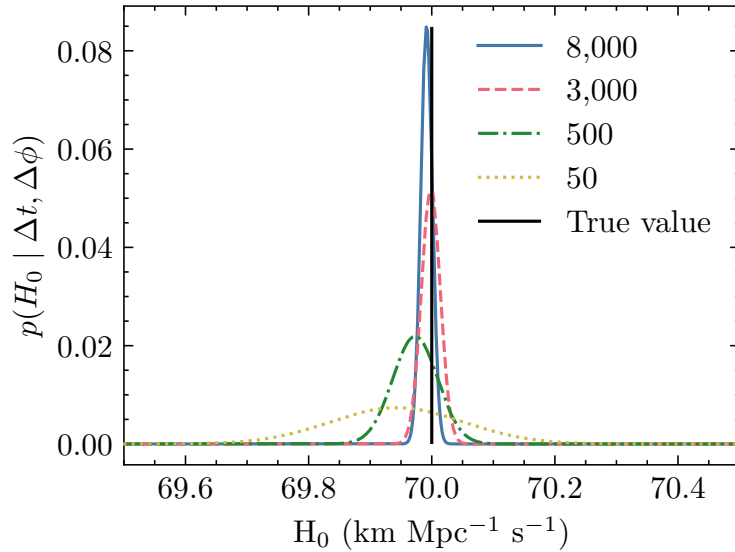
To assess the NRE statistical consistency, we conduct a coverage diagnostic on 10,000 noisy examples from the test set. This diagnostic verifies if, for a given fraction of the examples, the true value falls within the equivalent credible interval of the estimated posterior. For example, for 68% of the examples, the true value must fall within the 68% probability interval of the NRE posterior for the estimator to be consistent. To construct the credible intervals, we consider that each truth is part of its highest posterior density (HPD) interval, i.e. the interval where the posterior density is equal or above the one at the truth. Furthermore, the examples selected for the diagnostic have  $65 \text{ km s}^{-1} \text{ Mpc}^{-1} < H_0 < 75 \text{ km s}^{-1} \text{ Mpc}^{-1}$ , so



**Fig. 5.3.** Twelve  $H_0$  inferences on examples from the test set with zero noise realization. The two first rows show results for six different doubles, and the two last rows, for six different quads. The true  $H_0$  value is also different on each plot. The blue dashed line indicates the true posterior distribution (computed with nested sampling and the true likelihood), the grey dotted vertical line represents the true value, and the red solid line is the NRE posterior distribution. For doubles, the neural ratio estimates follow closely the true posteriors. For the quads, the difference is more pronounced, but the two distributions still agree well with each other.



**Fig. 5.4.** Coverage diagnostic of the NRE. A perfectly consistent distribution would fall on the dashed line. An underconfident distribution, i.e. which overestimates its uncertainty, would lay on the top-left area. An overconfident distribution, which underestimates its uncertainty, would be in the bottom-right region. The NRE coverage, represented by the orange solid line, indicates a minor underconfident behaviour.



**Fig. 5.5.** Population inferences of  $H_0$  with the NRE. The blue solid line, the pink dashed line, the green dashed-dotted line, and the yellow dotted line represent populations of 8,000, 3,000, 500 and 50 lensing systems, respectively. The true value  $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$  is indicated by the vertical black solid line. It falls inside the  $2\sigma$  interval for all populations.

that the posterior tails are fully encompassed by the prior on which the NRE was trained. Figure 5.4 displays the result of the coverage diagnostic. A perfectly consistent distribution would follow the black dashed line. The NRE shows a slightly underconfident behaviour,



which is not ideal, but preferable to an overconfident one. This mostly means that the NRE would need more observations to achieve the same precision as a calibrated estimator.

We also perform a population inference of  $H_0$  to verify that the NRE is unbiased. We simulate noisy data from multiple lensing systems, doubles as well as quads, but always with  $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$ . The inference on a single system is independent from all the others. Therefore, the joint posterior is the product of all individual posteriors. Figure 5.5 reports the results for populations of 8,000, 3,000, 500 and 50 lensing systems, respectively as the blue solid line, the pink dashed line, the green dashed-dotted line, and the yellow dotted line. The true value is represented by the vertical black solid line. The NRE appears unbiased because all posteriors enclose the truth in their  $2\sigma$  interval.

## 5.6. Discussion and Conclusion

Previous works have also explored how machine learning can improve the determination of  $H_0$  with time delay cosmography. Park et al. (2021) performed a joint inference on a large set of lensing systems with simulated time delay measurements and lensing parameter posterior distributions modeled by a BNN. Hence, as opposed of being used directly in the  $H_0$  inference, the trained BNN was rather a way to quickly execute the lens modeling and to obtain the parameter distributions, not only point estimates. From the BNN parameter distributions, it is straight-forward to compute the Fermat potential at the image positions. This approach is therefore very complementary to what is proposed here to estimate the  $H_0$  posterior with a NRE from modeled Fermat potentials and time delay measurements.

Moreover, one of the main advantages of a simulation-based approach such as the NRE over traditional maximum-likelihood methods is that it implicitly marginalizes over nuisance parameters. This provides a major simplification for inference problems in low dimension with a large number of nuisance parameters, such as the inference of  $H_0$  with time delay cosmography. For example, by replacing the Fermat potentials with the lensing observation in the inputs, the NRE could marginalize over the lensing parameters, removing the need for explicit lens modeling, and thus greatly accelerating the analysis. The training should still enforce the extraction of features related to the Fermat potentials from the lensing observation, considering that they are necessary information to discern the  $H_0$  contribution to the time delays. Nonetheless, it could be an interesting avenue for future work.

Another advantage of implicit-likelihood methods, equally important, is that one does not need to make any assumption about the form of the posterior to perform the inference. The complexity of the posterior is only limited by the simulations, which can include arbitrary complex environment, noise, selection effects, and systematics. Their realism is of the utmost importance, because the NRE learns the distribution ratio from them. By contrast,

traditional explicit-likelihood methods require an analytical form for both the prior and the likelihood to compute the posterior distribution. These often imply simplistic priors, and simplifying assumptions about the parametrization of the model, which can lead to various biases in the inference.

A notable source of bias given a too simplistic prior is the mass sheet degeneracy (Falco et al., 1985). It arises from a uniform mass density, whose lensing effect impacts the time delays, but leaves invariant the other observables. Thus, its influence can be confused with that of  $H_0$ , leading to a biased inference. In practice, the mass-sheet degeneracy is caused by line-of-sight perturbers whose effect is small enough for them to be approximated as an external convergence in the main deflector plane. Different mass profiles with equally good fit can also yield such a degeneracy. In this paper, we do not consider explicitly the mass sheet degeneracy. However, we chose the noise distributions so that the uncertainty on  $H_0$  could reach 8% frequently, which is the error budget estimated by Birrer, Simon & Treu, Tommaso (2021) when taking the mass sheet degeneracy into account.

In this work, we perform time delay cosmography on multiple strong lensing systems with an amortized Neural Ratio Estimator. We consider a variable point source and a SIE mass profile to simulate true Fermat potentials and time delays between the images. To keep the likelihood tractable, we emulate with Gaussian noise the analysis products of lens modeling and of time delay measurements. These are fed to the NRE, which has a Set Transformer architecture. The NRE posteriors are in agreement with the ones obtained with the traditional method. A coverage diagnostic indicates that the estimator slightly overestimates the uncertainties, and a population inference suggests that it is unbiased. The remaining improvements include more realistic simulations, the use of the raw data (the lensing observation and the light curves), and the incorporation of the mass sheet degeneracy in this framework.

# Chapitre 6

---

## Conclusion

La constante de Hubble a joué un rôle primordial dans l'histoire de la cosmologie moderne. Ses premières estimations, bien qu'erronées, correspondent à la découverte de l'expansion de l'Univers. Ce faisant, elles ont soutenu l'idée naissante d'un Univers dynamique et la conception d'une métrique pour décrire un tel espace-temps. Dans le but d'améliorer la mesure de  $H_0$ , les méthodes d'évaluation des distances astrophysiques se sont raffinées, notamment celle employant les supernovae de type Ia. Cette technique a mené à la découverte de l'accélération de l'expansion de l'Univers, et du même fait, à l'avènement du modèle cosmologique encore accepté à ce jour,  $\Lambda$ CDM. Ironiquement, la constante de Hubble incarne le doute qui plane sur ce modèle qu'elle a aidé à établir. En effet, la tension entre les mesures de  $H_0$  dépendantes et indépendantes de  $\Lambda$ CDM, en plus d'autres failles, suggère fortement que celui-ci est incomplet. Ce désaccord alimente donc activement nombre de recherches en cosmologie.

Parmi ces champs d'expertise stimulés par la crise de la cosmologie, on compte les lentilles gravitationnelles. Une image déformée et amplifiée résulte de ce phénomène, puisqu'elle suit la courbure de l'espace-temps infligée par une distribution de masse. L'amplification permet une résolution accrue des sources lumineuses lointaines et faibles. Quant à la déformation, elle renseigne sur la répartition de la masse du déflecteur. Surout, le temps de voyage de la lumière subissant un effet de lentille gravitationnelle dépend de la constante de Hubble via une combinaison de distances angulaires de diamètre. Cette relation pourrait jeter un nouvel éclairage sur le conflit concernant  $H_0$ , puisque les lentilles gravitationnelles sont indépendantes de  $\Lambda$ CDM et de l'échelle cosmique des distances, laquelle intervient dans la vaste majorité des mesures directes. L'inférence de la constante de Hubble à l'aide des lentilles gravitationnelles comporte tout de même certains défis dont une optimisation coûteuse en temps et en ressources de calcul.

Toutefois, les simulations donnent accès à un large volume de données. Cette condition est favorable à l'application de l'apprentissage automatique, lequel accélère significativement les

tâches d’optimisation. Il offre aussi des outils d’inférence comme les estimateurs neuronaux de ratio. Ces derniers apprennent à distinguer des exemples tirés de la distribution jointe entre paramètres d’intérêt et observables ainsi que du produit entre les distributions marginales de ceux-ci. La fonction ainsi obtenue équivaut au ratio entre la fonction de vraisemblance et l’évidence, ce qui correspond à la distribution à posteriori lorsque multiplié par celle a priori. Ces estimateurs se démarquent par leur généralisation, leur flexibilité et par leur marginalisation implicite des paramètres de nuisance. Par ailleurs, l’apprentissage automatique fournit des architectures qui sont invariantes sous la permutation de leurs entrées et qui en acceptent un nombre variable. C’est le cas du Set Transformer, ce qui le rend particulièrement commode pour l’inférence de  $H_0$  car les données utilisées, soit les temps de voyage et les potentiels de Fermat relatifs entre les images, forment un ensemble.

Dans ce mémoire, nous avons appliqué un estimateur neuronal de ratio à l’inférence de la constante de Hubble à partir de simulations de lentilles gravitationnelles fortes. Notre but consistait à montrer la faisabilité de cette technique. Nous avons donc misé sur la simplicité plutôt que le réalisme. Un ellipsoïde isotherme singulier caractérisait la distribution de masse surfacique du déflecteur, tandis que du bruit gaussien imitait les incertitudes générées par la modélisation du système de lentille et par la mesure des temps de voyage relatifs. Les distributions estimées par le NRE sont similaires à celles obtenues par une méthode conventionnelle, le *nested sampling*. L’estimateur ne présente pas de biais, mais il surestime légèrement son incertitude.

Afin d’appliquer notre méthode à de véritables observations, le travail se poursuit sur plusieurs aspects. D’abord, le réalisme des simulations d’entraînement doit augmenter pour s’assurer que l’estimateur neuronal de ratio apprend une distribution près de la vérité. Le réseau de neurones devra d’ailleurs conserver sa performance sur ces données beaucoup plus complexes. Une autre avenue à explorer serait l’utilisation des données brutes, soit les courbes de lumière des images et les observations de lentilles gravitationnelles, au lieu des produits d’analyse que sont les mesures des délais et les modélisations des potentiels de Fermat. De cette façon, l’inférence exploiterait pleinement la marginalisation implicite de l’estimateur neuronal de ratio. De plus, le temps requis par l’analyse complète s’en verrait encore réduit. Par ailleurs, un traitement rigoureux de la dégénérescence de la feuille de masse s’impose. Une entrée supplémentaire pourrait s’ajouter au réseau de neurones, laquelle correspondrait à une donnée qui brise une certaine version de la dégénérescence, telle que la vitesse de dispersion stellaire. Sinon, entraîner l’estimateur neuronal de ratio sur des simulations comportant des dégénérescences permettrait au moins de quantifier l’incertitude qu’elles engendrent.

Ce travail présente un moyen d’améliorer la vitesse et l’exactitude de la méthode de mesure de  $H_0$  basée sur les lentilles gravitationnelles. Ces contributions servent ainsi les intérêts

actuels de la cosmologie mentionnés au chapitre 1. De plus, ce mémoire est un exemple supplémentaire de la pertinence de combiner l'apprentissage automatique et l'inférence par simulations. L'architecture choisie pour le réseau de neurones artificiels, le Set Transformer, s'est avéré efficace pour traiter des ensembles, soit un type de données fréquent en astrophysique. Ainsi, ce travail s'inscrit aussi dans l'intégration de l'apprentissage automatique aux recherches en astrophysique. Cette technologie, bien que reconnue comme puissante, doit faire ses preuves pour gagner la confiance des scientifiques.



## Références bibliographiques

---

- Adam, A., Perreault-Levasseur, L., & Hezaveh, Y. 2022, Pixelated Reconstruction of Gravitational Lenses using Recurrent Inference Machines, arXiv, doi: 10.48550/ARXIV.2207.01073
- Arnett, W. D. 1969, *Astrophysics and Space Science*, 5, 180, doi: 10.1007/BF00650291
- Ba, J. L., Kiros, J. R., & Hinton, G. E. 2016, Layer Normalization, arXiv, doi: 10.48550/ARXIV.1607.06450
- Bartelmann, M., & Schneider, P. 2001, *Physics Reports*, 340, 291, doi: [https://doi.org/10.1016/S0370-1573\(00\)00082-X](https://doi.org/10.1016/S0370-1573(00)00082-X)
- Beaulieu, J.-P., Bennett, D. P., Fouqué, P., et al. 2006, *Nature*, 439, 437, doi: 10.1038/nature04441
- Behr, A. 1951, *Astronomische Nachrichten*, 279, 97, doi: <https://doi.org/10.1002/asna.19512790301>
- Birrer, Simon, & Treu, Tommaso. 2021, *Astronomy & Astrophysics*, 649, A61, doi: 10.1051/0004-6361/202039179
- Bishop, C. M. 2006, *Pattern Recognition and Machine Learning* (Springer)
- Bond, I. A., Udalski, A., Jaroszyński, M., et al. 2004, *The Astrophysical Journal*, 606, L155, doi: 10.1086/420928
- Bonvin, V., Courbin, F., Suyu, S. H., et al. 2016, *Monthly Notices of the Royal Astronomical Society*, 465, 4914, doi: 10.1093/mnras/stw3006
- Choromanska, A., Henaff, M., Mathieu, M., Ben Arous, G., & LeCun, Y. 2015, in *Proceedings of Machine Learning Research*, Vol. 38, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, ed. G. Lebanon & S. V. N. Vishwanathan (San Diego, California, USA: PMLR), 192–204. <https://proceedings.mlr.press/v38/choromanska15.html>
- Chwolson, O. 1924, *Astronomische Nachrichten*, 221, 329, doi: <https://doi.org/10.1002/asna.19242212003>
- Colgate, S. A., & McKee, C. 1972, in *Stellar Evolution*, ed. H. Y. Chiu & A. Muriel, 307
- Collaboration Planck, Ade, P. A. R., Aghanim, N., et al. 2014, *Astronomy & Astrophysics*, 571, A16, doi: 10.1051/0004-6361/201321591

- Collaboration Planck, Ade, P. A. R., Aghanim, N., et al. 2016, *Astronomy & Astrophysics*, 594, A13, doi: 10.1051/0004-6361/201525830
- Collaboration Planck, Aghanim, N., Akrami, Y., et al. 2020, *Astronomy & Astrophysics*, 641, A6, doi: 10.1051/0004-6361/201833910
- Cranmer, K., Brehmer, J., & Louppe, G. 2020, *Proceedings of the National Academy of Sciences*, 117, 30055, doi: 10.1073/pnas.1912789117
- Cranmer, K., Pavez, J., & Louppe, G. 2015, arXiv e-prints, arXiv:1506.02169. <https://arxiv.org/abs/1506.02169>
- Di Valentino, E., Mena, O., Pan, S., et al. 2021, *Classical and Quantum Gravity*, 38, 153001, doi: 10.1088/1361-6382/ac086d
- Duchi, J., Hazan, E., & Singer, Y. 2011, *Journal of machine learning research*, 12
- Eddington, A. S. 1919, *The Observatory*, 42, 119
- Eddington, A. S. 1923, *The mathematical theory of relativity* (The University Press)
- Einstein, A. 1915, *Sitzungsber. Preuss. Akad. Wiss. Berlin (Math. Phys. )*, 1915, 844
- . 1936, *Science*, 84, 506. <http://www.jstor.org/stable/1663250>
- Falco, E. E., Gorenstein, M. V., & Shapiro, I. I. 1985, *Astrophysical Journal Letters*, 289, L1, doi: 10.1086/184422
- Flores, R. A., & Primack, J. R. 1994, *Astrophysical Journal Letters*, 427, L1, doi: 10.1086/187350
- Freedman, W. L., Madore, B. F., Scowcroft, V., et al. 2012, *The Astrophysical Journal*, 758, 24, doi: 10.1088/0004-637X/758/1/24
- Freedman, W. L., Madore, B. F., Gibson, B. K., et al. 2001, *The Astrophysical Journal*, 553, 47, doi: 10.1086/320638
- Friedmann, A. 1922, *Zeitschrift für Physik*, 10, 377, doi: 10.1007/BF01332580
- . 1924, *Zeitschrift für Physik*, 21, 326, doi: 10.1007/BF01328280
- Fukugita, M., Hogan, C. J., & Peebles, P. J. E. 1993, *Nature*, 366, 309, doi: 10.1038/366309a0
- Glorot, X., & Bengio, Y. 2010, in *Proceedings of Machine Learning Research*, Vol. 9, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, ed. Y. W. Teh & M. Titterton (Chia Laguna Resort, Sardinia, Italy: PMLR), 249–256. <https://proceedings.mlr.press/v9/glorot10a.html>
- Goldstein, D. A., & Nugent, P. E. 2016, *The Astrophysical Journal Letters*, 834, L5, doi: 10.3847/2041-8213/834/1/L5
- Goodfellow, I., Bengio, Y., & Courville, A. 2016, *Deep Learning* (MIT Press)
- Goodfellow, I. J., Vinyals, O., & Saxe, A. M. 2014, *Qualitatively characterizing neural network optimization problems*, arXiv, doi: 10.48550/ARXIV.1412.6544
- Gould, A., & Loeb, A. 1992, *The Astrophysical Journal*, 396, 104
- Haas, A. E. 1938, *Science*, 87, 195, doi: 10.1126/science.87.2252.195



- Hamuy, M., Maza, J., Phillips, M., et al. 1993, *Astronomical Journal*, 106, 2392, doi: 10.1086/116811
- Handley, W. J., Hobson, M. P., & Lasenby, A. N. 2015a, *Monthly Notices of the Royal Astronomical Society: Letters*, 450, L61, doi: 10.1093/mnrasl/slv047
- . 2015b, *Monthly Notices of the Royal Astronomical Society*, 453, 4385, doi: 10.1093/mnras/stv1911
- Hermans, J., Begy, V., & Louppe, G. 2020, in *Proceedings of Machine Learning Research*, Vol. 119, *Proceedings of the 37th International Conference on Machine Learning*, ed. H. D. III & A. Singh (PMLR), 4239–4248. <https://proceedings.mlr.press/v119/hermans20a.html>
- Hezaveh, Y. D., Levasseur, L. P., & Marshall, P. J. 2017, *Nature*, 548, 555, doi: 10.1038/nature23463
- Hinton, G., Srivastava, N., & Swersky, K. 2012, Cited on, 14, 2
- Hornik, K. 1991, *Neural Networks*, 4, 251, doi: [https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T)
- Hoyle, F., & Fowler, W. A. 1960, *Astrophysical Journal*, 132, 565, doi: 10.1086/146963
- Hubble, E. 1929, *Proceedings of the National Academy of Sciences*, 15, 168, doi: 10.1073/pnas.15.3.168
- Jacoby, G. H., Walker, A. R., & Ciardullo, R. 1990, in *Bulletin of the American Astronomical Society*, Vol. 22, 866
- Jha, S., Garnavich, P. M., Kirshner, R. P., et al. 1999, *The Astrophysical Journal Supplement Series*, 125, 73, doi: 10.1086/313275
- Kingma, D. P., & Ba, J. 2014, *Adam: A Method for Stochastic Optimization*, arXiv, doi: 10.48550/ARXIV.1412.6980
- Klypin, A., Kravtsov, A. V., Valenzuela, O., & Prada, F. 1999, *The Astrophysical Journal*, 522, 82, doi: 10.1086/307643
- Kochanek, C. S. 2006, *Strong Gravitational Lensing* (Berlin, Heidelberg: Springer Berlin Heidelberg), 91–268, doi: 10.1007/978-3-540-30310-7\_2
- Kormann, R., Schneider, P., & Bartelmann, M. 1994, *Astronomy & Astrophysics*, 284, 285
- Kowal, C. T. 1968, *Astronomical Journal*, 73, 1021, doi: 10.1086/110763
- Lee, J., Lee, Y., Kim, J., et al. 2019, in *Proceedings of Machine Learning Research*, Vol. 97, *Proceedings of the 36th International Conference on Machine Learning*, ed. K. Chaudhuri & R. Salakhutdinov (PMLR), 3744–3753. <https://proceedings.mlr.press/v97/lee19d.html>
- Legin, R., Hezaveh, Y., Levasseur, L. P., & Wandelt, B. 2021, *Simulation-Based Inference of Strong Gravitational Lensing Parameters*, arXiv, doi: 10.48550/ARXIV.2112.05278
- Legin, R., Hezaveh, Y., Perreault-Levasseur, L., & Wandelt, B. 2022, *A Framework for Obtaining Accurate Posteriors of Strong Gravitational Lensing Parameters with Flexible*

- Priors and Implicit Likelihoods using Density Estimation, arXiv, doi: 10.48550/ARXIV.2212.00044
- Lemaître, G. 1927, *Annales de la Société Scientifique de Bruxelles*, A47, 49
- Levasseur, L. P., Hezaveh, Y. D., & Wechsler, R. H. 2017, *The Astrophysical Journal*, 850, L7, doi: 10.3847/2041-8213/aa9704
- Lundmark, K. 1924, *Monthly Notices of the Royal Astronomical Society*, 84, 747, doi: 10.1093/mnras/84.9.747
- . 1925, *MNRSA*, 85, 865, doi: 10.1093/mnras/85.8.865
- Mao, S., & Paczynski, B. 1991, *The Astrophysical journal*, 374, L37
- Massey, R., Kitching, T., & Richard, J. 2010, *Reports on Progress in Physics*, 73, 086901, doi: 10.1088/0034-4885/73/8/086901
- Meneghetti, M. 2013, *Introduction to Gravitational Lensing*
- Milne, E. A. 1933, *Monthly Notices of the Royal Astronomical Society*, 94, 3, doi: 10.1093/mnras/94.1.3
- Moore, B. 1994, *Nature*, 370, 629, doi: 10.1038/370629a0
- Moore, B., Ghigna, S., Governato, F., et al. 1999, *The Astrophysical Journal*, 524, L19, doi: 10.1086/312287
- Morningstar, W. R., Levasseur, L. P., Hezaveh, Y. D., et al. 2019, *The Astrophysical Journal*, 883, 14, doi: 10.3847/1538-4357/ab35d7
- Mould, J. R., Huchra, J. P., Freedman, W. L., et al. 2000, *The Astrophysical Journal*, 529, 786, doi: 10.1086/308304
- Oguri, M., & Marshall, P. J. 2010, *Monthly Notices of the Royal Astronomical Society*, no, doi: 10.1111/j.1365-2966.2010.16639.x
- Park, J. W., Wagner-Carena, S., Birrer, S., et al. 2021, *The Astrophysical Journal*, 910, 39, doi: 10.3847/1538-4357/abdfc4
- Pearson, J., Li, N., & Dye, S. 2019, *Monthly Notices of the Royal Astronomical Society*, 488, 991, doi: 10.1093/mnras/stz1750
- Perlmutter, S., Pennypacker, C. R., Goldhaber, G., et al. 1995, *Astrophysical Journal Letters*, 440, L41, doi: 10.1086/187756
- Perlmutter, S., Aldering, G., Goldhaber, G., et al. 1999, *The Astrophysical Journal*, 517, 565, doi: 10.1086/307221
- Phillips, M. M. 1993, *Astrophysical Journal Letters*, 413, L105, doi: 10.1086/186970
- Polyak, B. 1964, *USSR Computational Mathematics and Mathematical Physics*, 4, 1, doi: [https://doi.org/10.1016/0041-5553\(64\)90137-5](https://doi.org/10.1016/0041-5553(64)90137-5)
- Refsdal, S. 1964, *Monthly Notices of the Royal Astronomical Society*, 128, 307, doi: 10.1093/mnras/128.4.307
- Riess, A. G., Casertano, S., Yuan, W., et al. 2021, *The Astrophysical Journal Letters*, 908, L6, doi: 10.3847/2041-8213/abdbaf

- Riess, A. G., Casertano, S., Yuan, W., Macri, L. M., & Scolnic, D. 2019, *The Astrophysical Journal*, 876, 85, doi: 10.3847/1538-4357/ab1422
- Riess, A. G., Press, W. H., & Kirshner, R. P. 1996, *The Astrophysical Journal*, 473, 88, doi: 10.1086/178129
- Riess, A. G., Filippenko, A. V., Challis, P., et al. 1998, *The Astronomical Journal*, 116, 1009, doi: 10.1086/300499
- Riess, A. G., Macri, L., Casertano, S., et al. 2011, *The Astrophysical Journal*, 730, 119, doi: 10.1088/0004-637X/730/2/119
- Riess, A. G., Macri, L. M., Hoffmann, S. L., et al. 2016, *The Astrophysical Journal*, 826, 56, doi: 10.3847/0004-637X/826/1/56
- Riess, A. G., Casertano, S., Yuan, W., et al. 2018, *The Astrophysical Journal*, 855, 136, doi: 10.3847/1538-4357/aaadb7
- Riess, A. G., Yuan, W., Macri, L. M., et al. 2022, *The Astrophysical Journal Letters*, 934, L7, doi: 10.3847/2041-8213/ac5c5b
- Robertson, H. P. 1933, *Rev. Mod. Phys.*, 5, 62, doi: 10.1103/RevModPhys.5.62
- Robertson, H. P. 1935, *Astrophysical Journal*, 82, 284, doi: 10.1086/143681
- . 1936a, *Astrophysical Journal*, 83, 187, doi: 10.1086/143716
- . 1936b, *Astrophysical Journal*, 83, 257, doi: 10.1086/143726
- Sandage, A., Tammann, G. A., Saha, A., et al. 2006, *The Astrophysical Journal*, 653, 843, doi: 10.1086/508853
- Schmidt, B. P., Suntzeff, N. B., Phillips, M. M., et al. 1998, *The Astrophysical Journal*, 507, 46, doi: 10.1086/306308
- Schneider, & Sluse. 2013, *Astronomy & Astrophysics*, 559, A37, doi: 10.1051/0004-6361/201321882
- . 2014, *Astronomy & Astrophysics*, 564, A103, doi: 10.1051/0004-6361/201322106
- Schneider, P. 2006, *Introduction to Gravitational Lensing and Cosmology* (Berlin, Heidelberg: Springer Berlin Heidelberg), 1–89, doi: 10.1007/978-3-540-30310-7\_1
- Schneider, P. 2019, *Astronomy & Astrophysics*, 624, A54, doi: 10.1051/0004-6361/201424881
- Schuldt, S., Suyu, S. H., Meinhardt, T., et al. 2021, *Astronomy & Astrophysics*, 646, A126, doi: 10.1051/0004-6361/202039574
- Schöneberg, N., Abellán, G. F., Sánchez, A. P., et al. 2022, *Physics Reports*, 984, 1, doi: <https://doi.org/10.1016/j.physrep.2022.07.001>
- Silberstein, L. 1924, *Nature*, 114, 347, doi: 10.1038/114347b0
- Slipher, V. M. 1917, *Proceedings of the American Philosophical Society*, 56, 403
- Spergel, D. N., Verde, L., Peiris, H. V., et al. 2003, *The Astrophysical Journal Supplement Series*, 148, 175, doi: 10.1086/377226

- Suyu, S. H., Chang, T.-C., Courbin, F., & Okumura, T. 2018, *Space Science Reviews*, 214, 91, doi: 10.1007/s11214-018-0524-3
- Tonry, J. 1991, *Sky and Telescope*, 82, 460
- Treu, T. 2010, *Annual Review of Astronomy and Astrophysics*, 48, 87, doi: 10.1146/annurev-astro-081309-130924
- Treu, T., Roberts-Borsani, G., Bradac, M., et al. 2022, *The Astrophysical Journal*, 935, 110, doi: 10.3847/1538-4357/ac8158
- Trimble, V. 1996, *Publications of the Astronomical Society of the Pacific*, 108, 1073, doi: 10.1086/133837
- Tully, R. B., & Fisher, J. R. 1977, *Astronomy and Astrophysics*, 54, 661
- Vaswani, A., Shazeer, N., Parmar, N., et al. 2017, in *Advances in Neural Information Processing Systems*, ed. I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett, Vol. 30 (Curran Associates, Inc.). <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>
- Walker, A. G. 1937, *Proceedings of the London Mathematical Society*, s2-42, 90, doi: <https://doi.org/10.1112/plms/s2-42.1.90>
- Walsh, D., Carswell, R. F., & Weymann, R. J. 1979, *Nature*, 279, 381, doi: 10.1038/279381a0
- Weyl, H. 1923, *General Relativity and Gravitation*, 41, 1661, doi: 10.1007/s10714-009-0826-6
- Zaheer, M., Kottur, S., Ravanbakhsh, S., et al. 2017, *Deep Sets*, arXiv, doi: 10.48550/ARXIV.1703.06114
- Zwicky, F. 1929, *Proceedings of the National Academy of Sciences*, 15, 773, doi: 10.1073/pnas.15.10.773
- . 1937, *Phys. Rev.*, 51, 290, doi: 10.1103/PhysRev.51.290

# Annexe A

---

## Neural network variable dimensions

**Table A.1.** Input sizes for each operation in the NRE

<b>Operation</b>	<b>Input sizes</b>
First multihead attention block	example set size $\times$ 2
Second multihead attention block	example set size $\times$ 384
Third multihead attention block	features : example set size $\times$ 384 learnable seed vector : example set size $\times$ 1 $\times$ 384
Concatenating $H_0$ and the redshifts	384
First linear layer	387
Second linear layer	768
Third linear layer	768
Ratio estimation (see Equation (5.2.12))	2



## Annexe B

---

### Congrès où l'étudiante a présenté ses résultats

#### Octobre numérique d'IVADO

**Médium** : Présentation orale

**Titre** : L'apprentissage automatique au service de la cosmologie en crise

**Lieu** : Montréal, QC

**Année** : 2022

**Auteur.e.s** : È. Campeau-Poirier, L. Perreault Levasseur, Y. Hezaveh

#### Boom! A Workshop on Explosive Transients with LSST

**Médium** : Présentation orale (en ligne)

**Titre** : Time delay cosmography with a neural ratio estimator

**Lieu** : Chicago, IL

**Année** : 2022

**Auteur.e.s** : È. Campeau-Poirier, L. Perreault Levasseur, Y. Hezaveh

#### Rencontre annuelle du CRAQ

**Médium** : Présentation orale

**Titre** : Time delay cosmography with a neural ratio estimator

**Lieu** : Magog, QC

**Année** : 2022

**Auteur.e.s** : È. Campeau-Poirier, L. Perreault Levasseur, Y. Hezaveh

#### Likelihood Free In Paris

**Médium** : Présentation orale

**Titre** : Time delay cosmography with a neural ratio estimator

**Lieu :** Paris, FR

**Année :** 2022

**Auteur.e.s :** È. Campeau-Poirier, L. Perreault Levasseur, Y. Hezaveh

## **Mon projet de recherche en 180 secondes d'IVADO**

**Médium :** Présentation orale

**Titre :** Mesurer le taux d'expansion de l'Univers à l'aide de l'apprentissage automatique

**Lieu :** Montréal, QC

**Année :** 2022

**Auteur.e.s :** È. Campeau-Poirier

## **Octobre numérique d'IVADO**

**Médium :** Présentation orale (en ligne)

**Titre :** Measuring the expansion rate of the Universe

**Lieu :** Montréal, QC

**Année :** 2021

**Auteur.e.s :** È. Campeau-Poirier, L. Perreault Levasseur, Y. Hezaveh