

AI for astronomical spectroscopy: current challenges and prospects for the WST

Nils Candebat

INAF- Observatory of Arcetri

MALSPEC : Machine Learning for Spectroscopy

Germano Sacco, Francesco Belfiore, Laura Magrini, Stefano Zibetti

nils.candebat@inaf.it

Painting in a watercolour style of a stunning view of a spiral galaxy, likely captured by a space telescope. The galaxy is centrally located and appears to be a bright, luminous core surrounded by swirling arms of gas and dust. The arms are composed of various shades of blue, white, and even hints of orange, indicating the presence of different elements and temperatures within the galaxy. The background is a deep, dark space, filled with numerous stars of varying brightness. The stars are scattered throughout the image, creating a sense of depth and vastness. The overall composition of the image highlights the intricate beauty and complexity of the galaxy, showcasing its spiral structure and



WST MALSPEC

Generated using DeepSeek Janus F

Neural Network in one slide

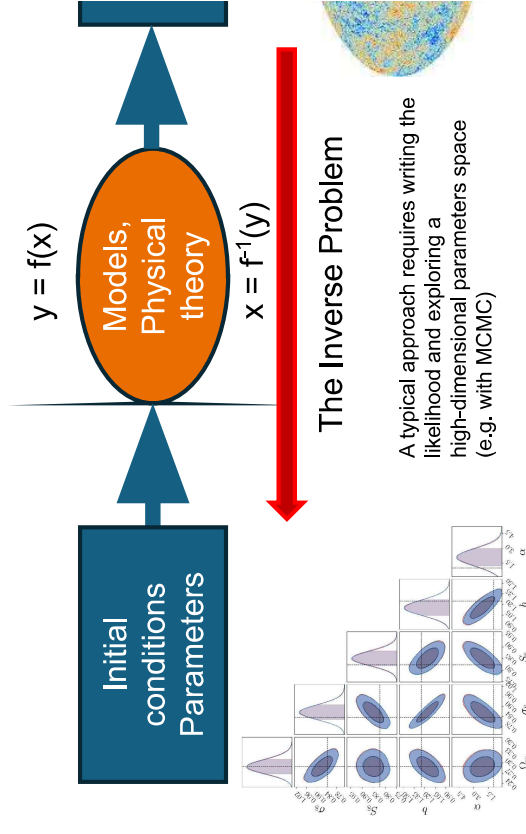
Universal Approximation Theorem:

Let $f : K \rightarrow \mathbb{R}$ be a continuous function, where K is a compact subset of \mathbb{R}^n . For any $\epsilon > 0$, there exists a feedforward neural network $\hat{f}(x)$ with a single hidden layer and a finite number of neurons such that:

$$|f(x) - \hat{f}(x)| < \epsilon \quad \forall x \in K$$

where $\hat{f}(x)$ is the output of the neural network. The network uses a non-constant, bounded, and continuous activation function σ .

Solving the inverse problem



cINN applied to stars

Drawing in an art deco style of a stunning view of the cosmos, showcasing a vast expanse of space filled with numerous stars and celestial bodies. The background is a deep blue, representing outer space, with stars of various sizes and colours scattered throughout. Some stars are bright white, while others are fainter and orange-brown, indicating different temperatures and compositions. In the centre of the image, there is a dense cluster of stars, with some forming a bright, glowing nucleus surrounded by smaller stars. This central region appears to be a galaxy or a nebula, with its intricate structure and glowing core being particularly prominent. The image also shows some nebulous formations, which are wispy clouds of gas and dust, adding to the complexity and beauty of the cosmic scene. These nebulous regions are scattered across the image, contributing to the overall sense of depth and vastness. @Hubble



Generated using DeepSeek Janus F

Dataset

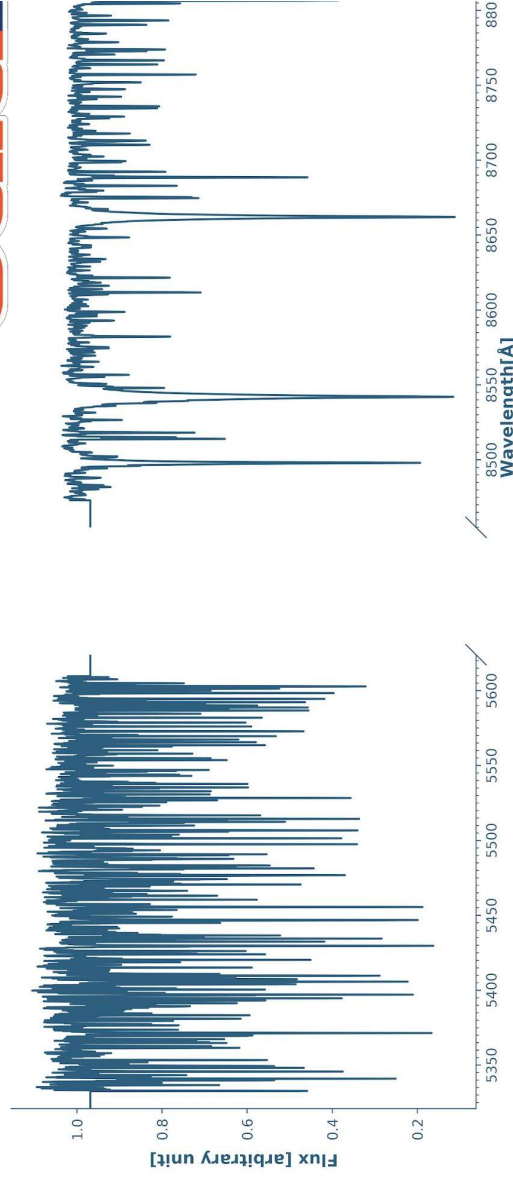
=> Spectra: GAIA-ESO Survey (GES), Randich et al. 2022

- 52,841 spectra from GIRAFFE@VLT
- spectral range 5330-5610Å & 8480-8980Å
- 17,000 pixels per spectrum



=> 9 Parameters (targets) derived from multiple classical pipelines Hourihane et al. (2022):

- Temperature effective - $TEFF$
- Surface gravity - $\log(g)$
- Metallicity abundance - $[Fe/H]$
- Aluminium abundance - $[Al/H]$
- Magnesium abundance - $[Mg/H]$
- Calcium abundance - $[Ca/H]$
- Nickel abundance - $[Ni/H]$
- Titanium abundance - $[Ti/H]$
- Silicon abundance - $[Si/H]$

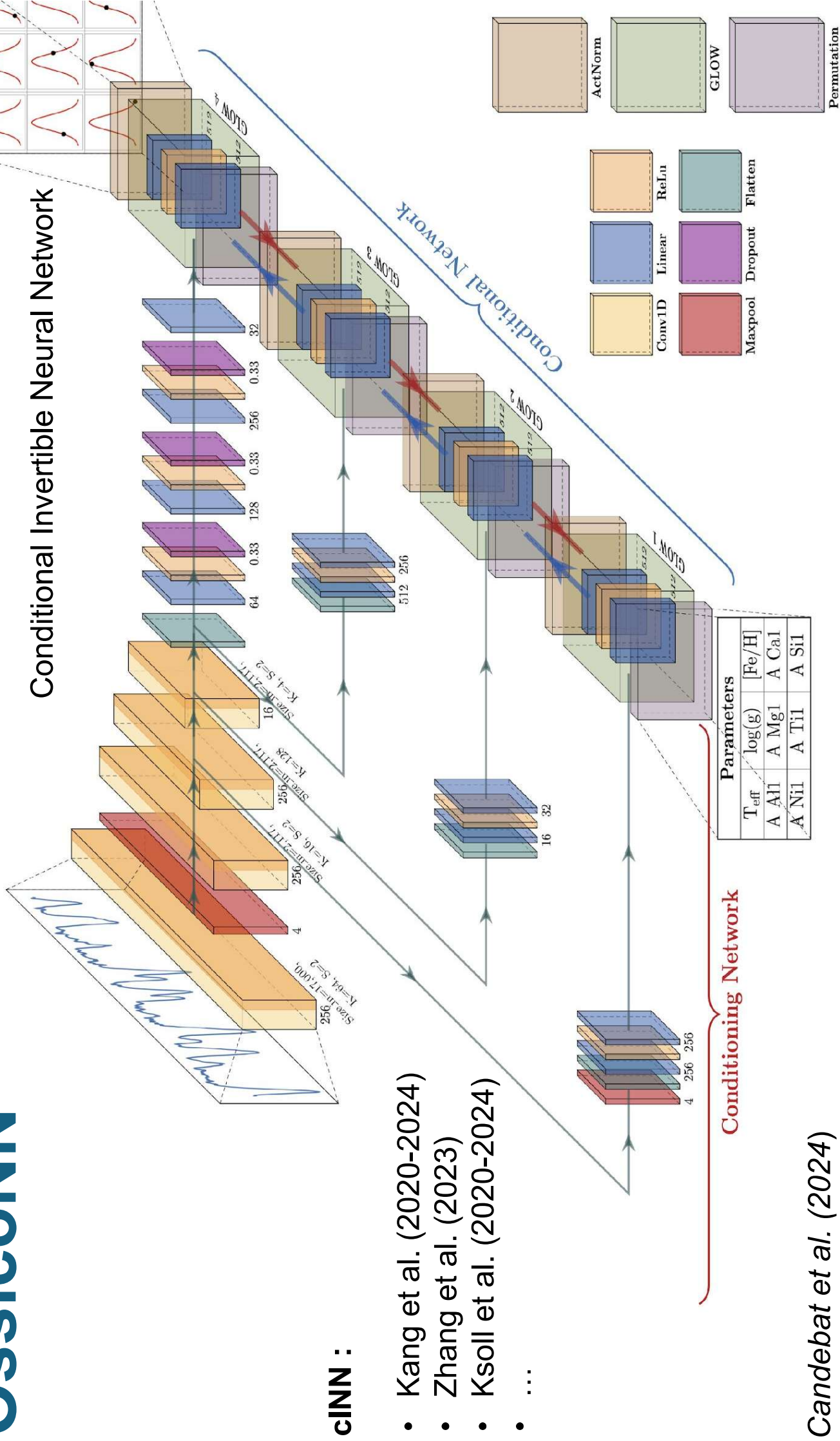


Dataset distribution:

- GIRAFFE HR10-HR21 : all stars [52,841 stars]
- Training-Test : high SNR(>25) and 9 params [6,963 stars]
- Noisy : low SNR (<25) and 9 params [2,613 stars]

Ossiconn

Conditional Invertible Neural Network



cINN :

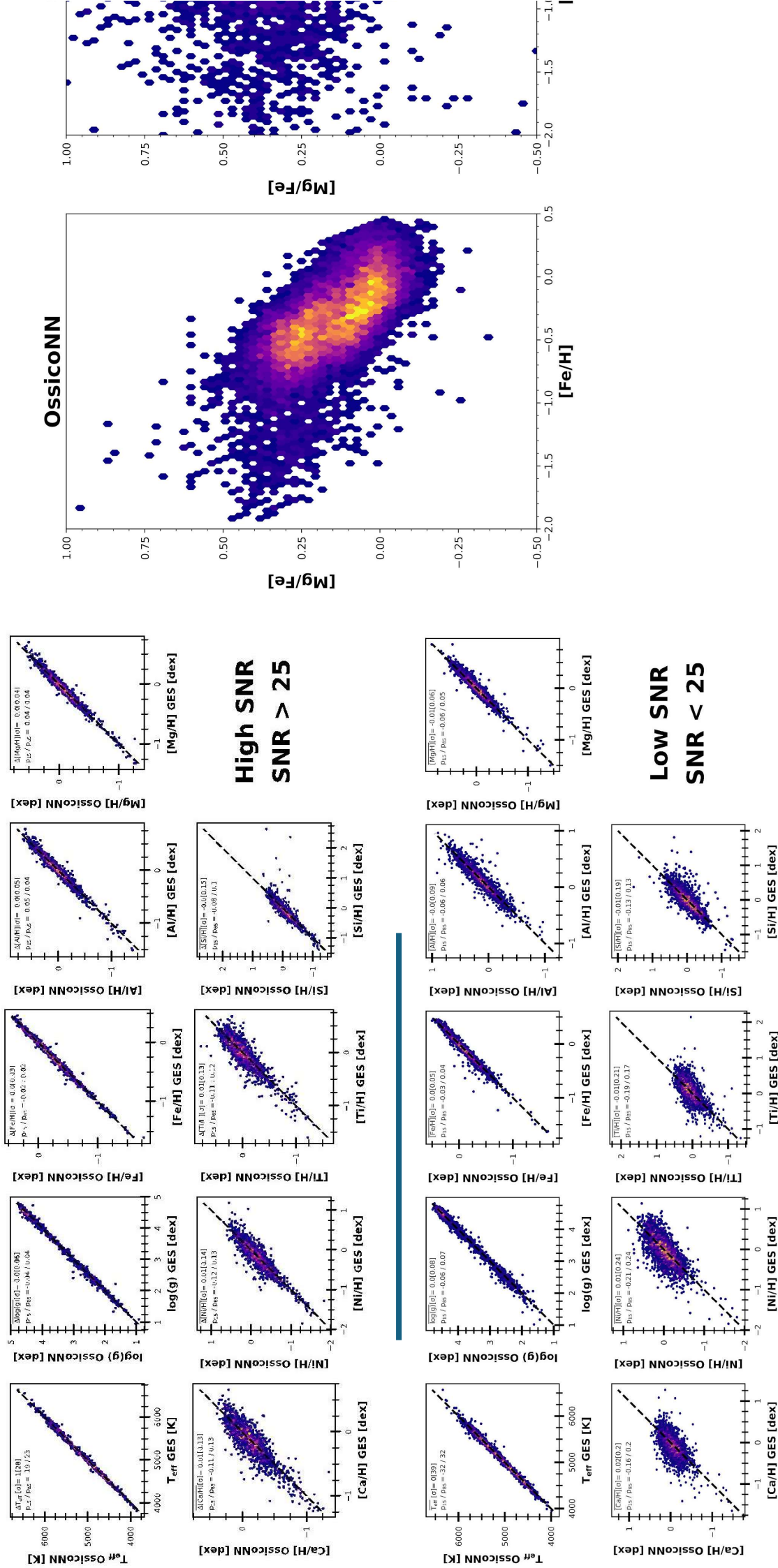
- Kang et al. (2020-2024)
- Zhang et al. (2023)
- Ksoil et al. (2020-2024)
- ...

Candebat et al. (2024)

AI & cINN

Applied to stars

Results

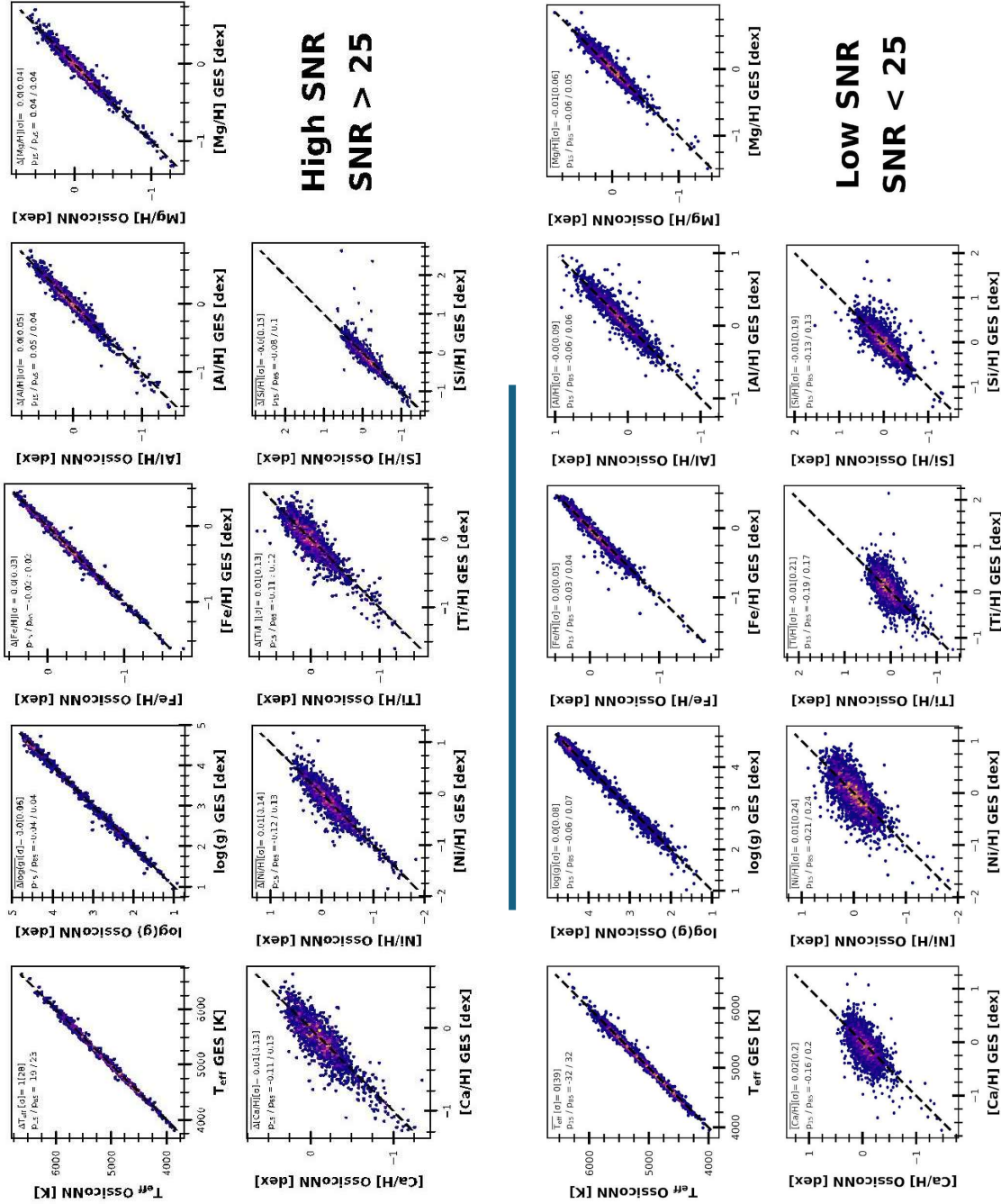


Candebat et al. (2024)

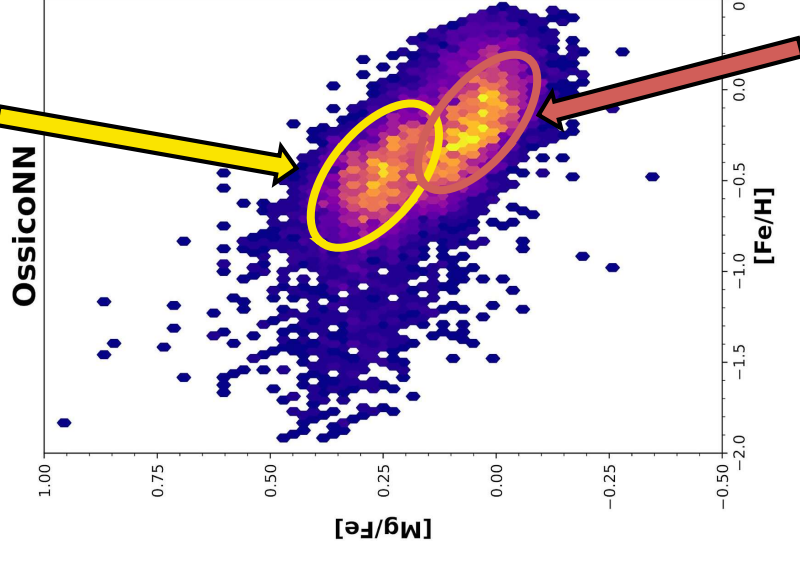
AI & cINN

Applied to stars

Results



Thick disk



Thin disk

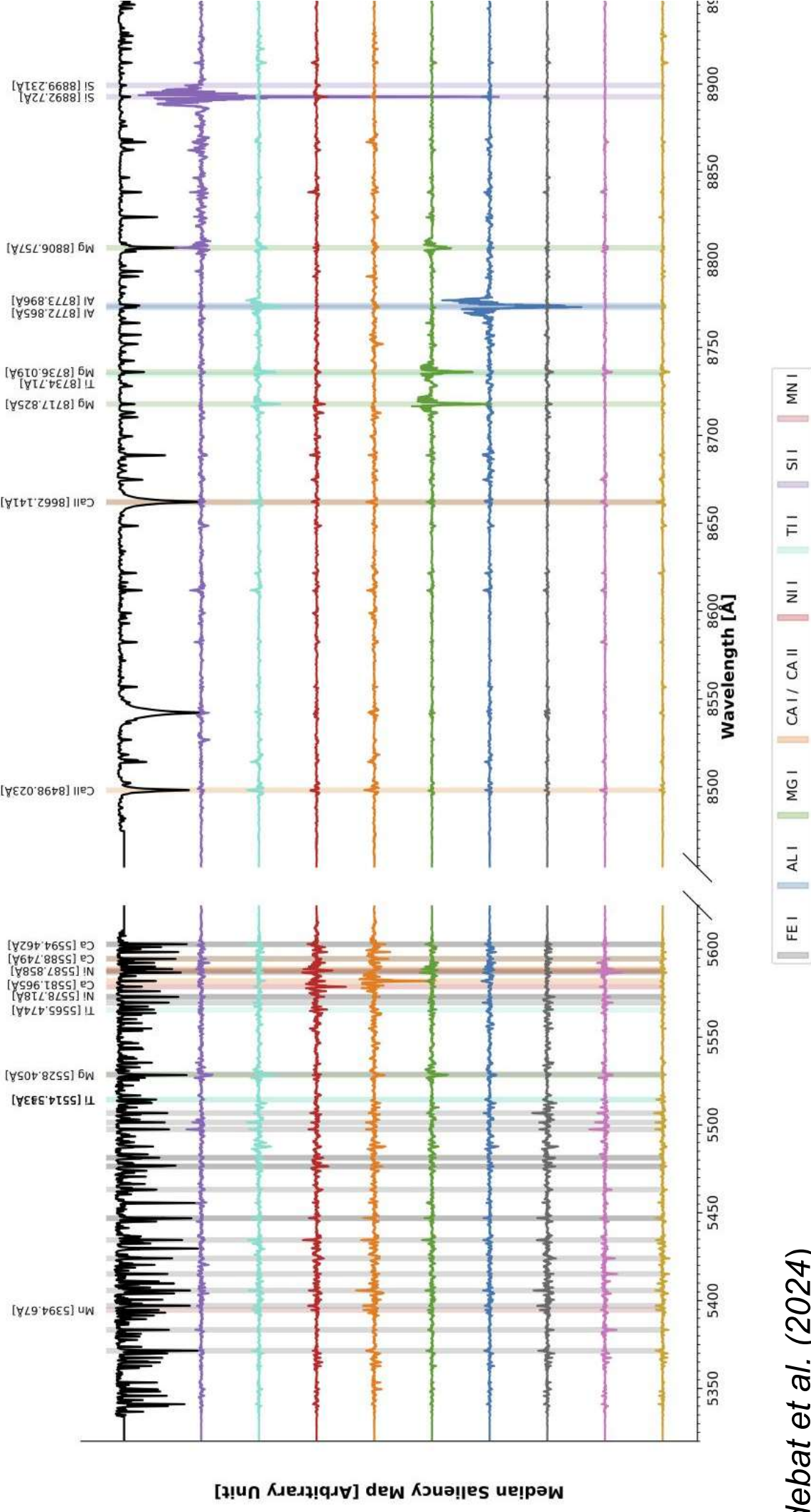
Candebat et al. (2024)

AI & cINN

Applied to stars

Alignment Validation: Saliency map

Saliency Map: technique to visualizing the input features that contribute the most to its output prediction

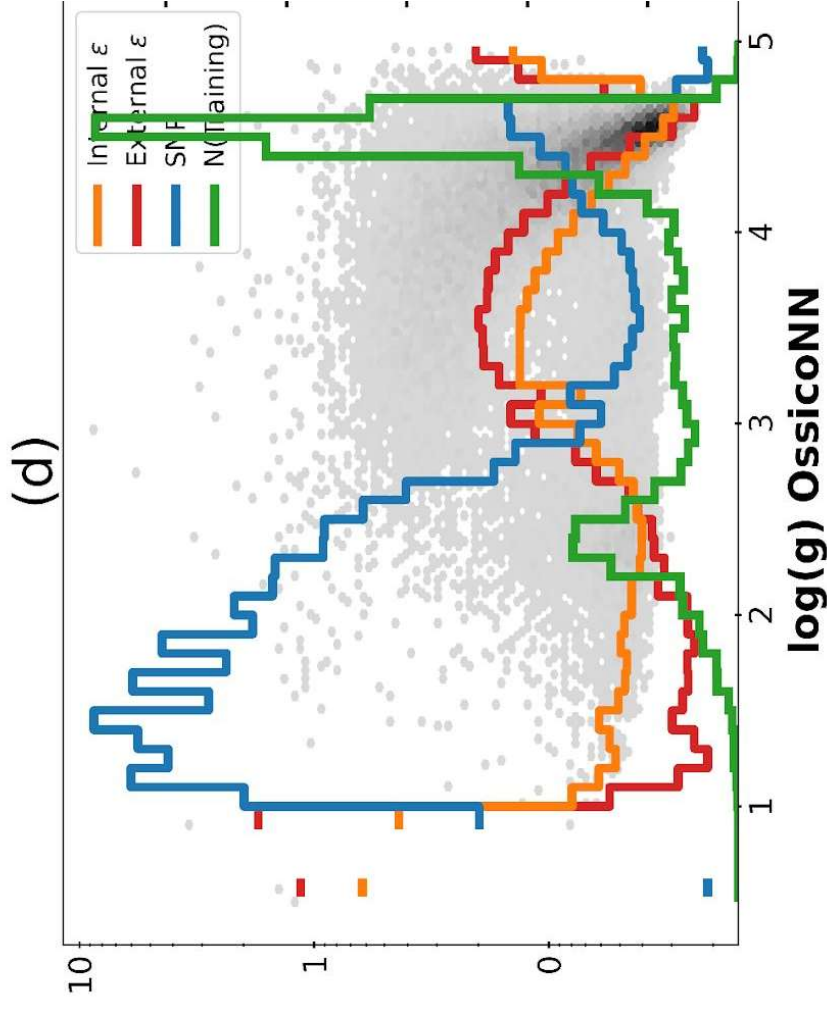
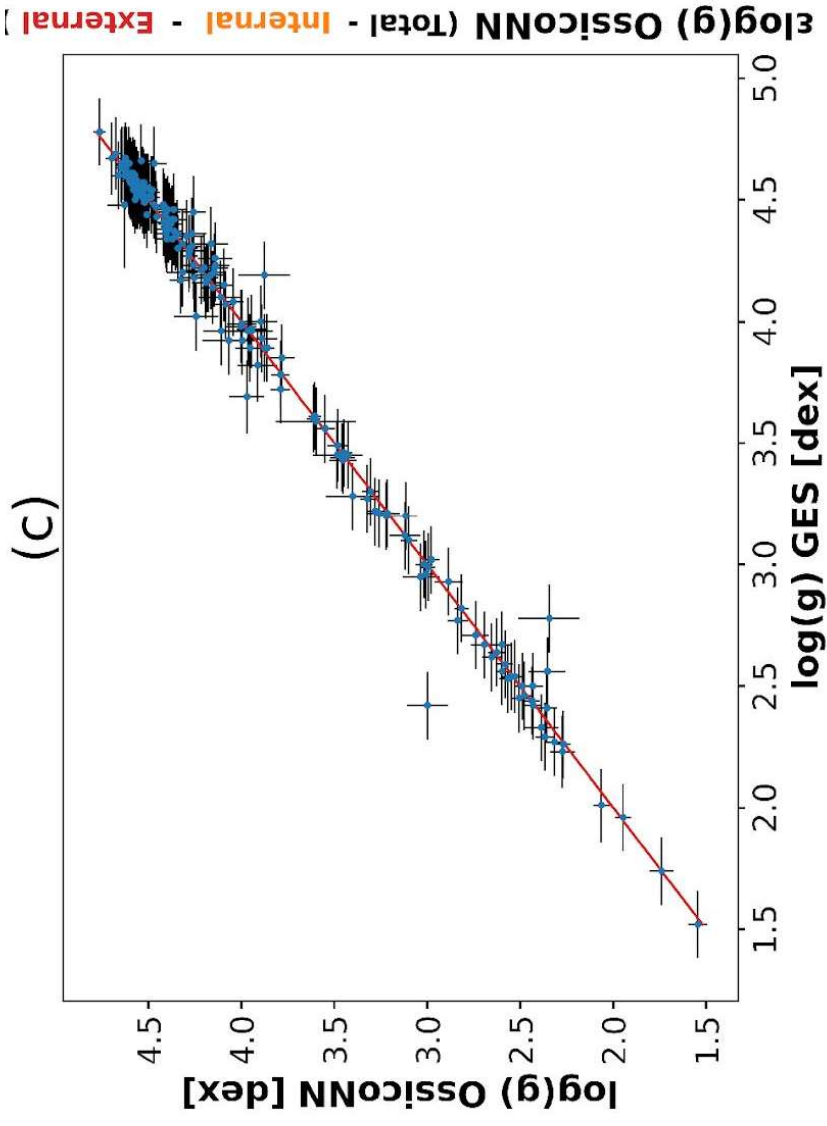


Candebat et al. (2024)

AI & cINN

Applied to stars

Uncertainty



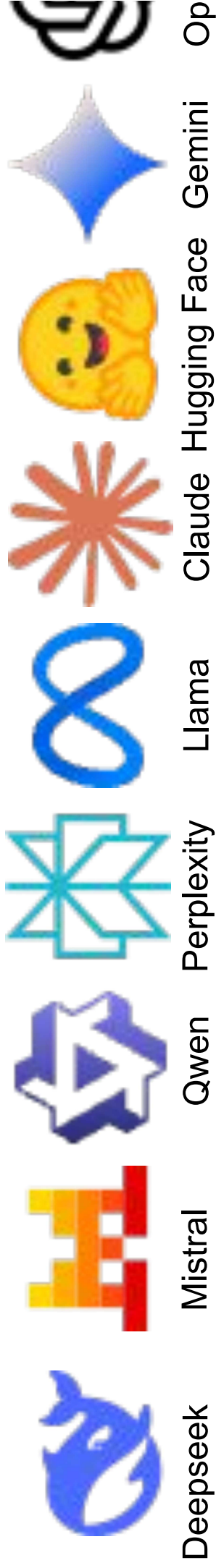
Large Language model Foundation Model

Drawing in an art deco style of a collection of galaxies captured in a deep space photograph. The galaxies are of varying shapes and sizes, with some appearing as spiral nebulae and others as elliptical or irregular shapes. The background is filled with numerous small stars, creating a dense and vibrant cosmic scene. The galaxies are illuminated in different colors, with some displaying a mix of blues, purples, and reds, indicating the presence of various gases and star formations. The overall composition suggests a rich and diverse galaxy cluster, showcasing the intricate beauty of the universe.

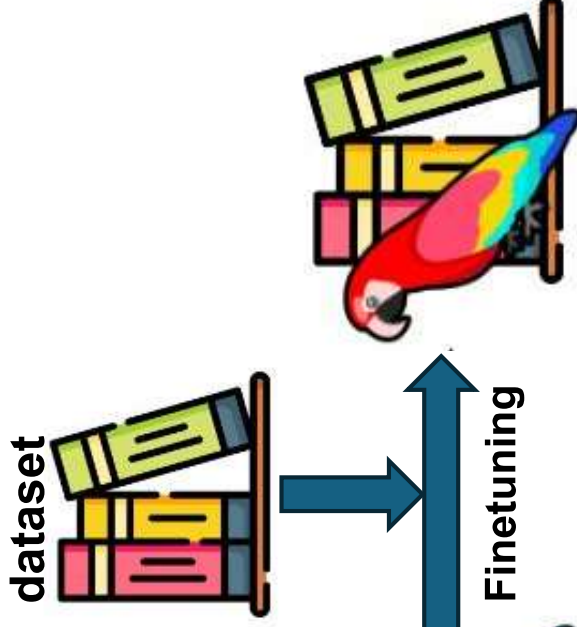


Generated using DeepSeek Janus F

Idea



Specific and small dataset



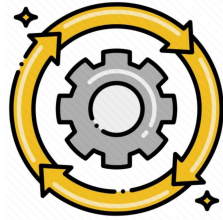
Foundation Model

Finetuned LLM

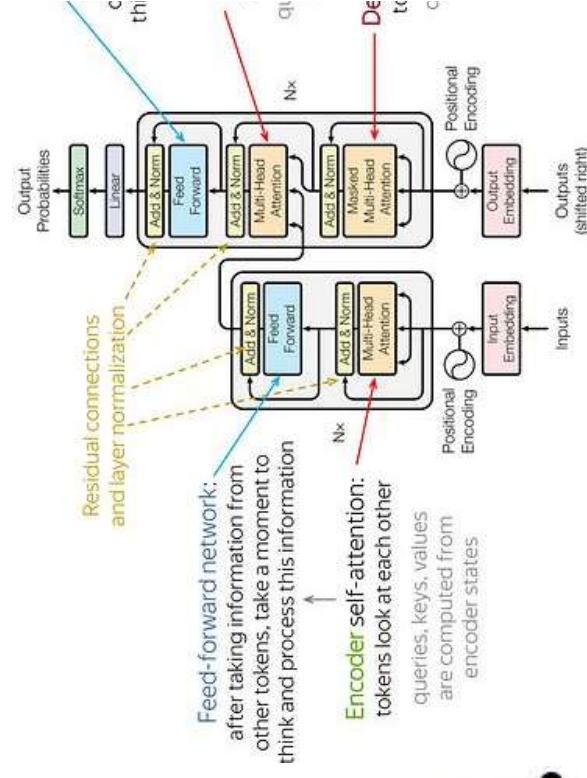
Tomaž Bratanič



Very large Dataset



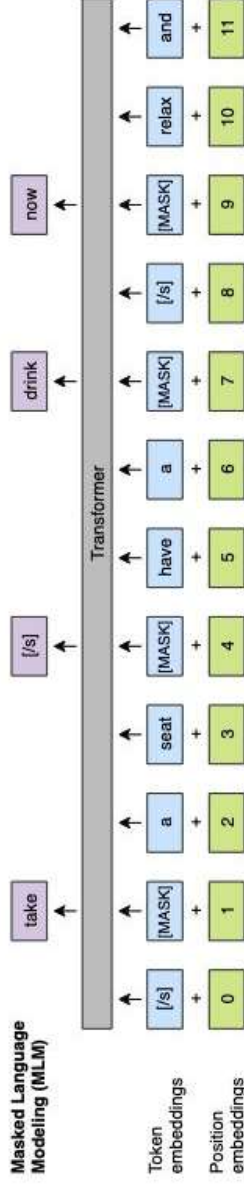
LLM Architecture



Transformers – Self Attention

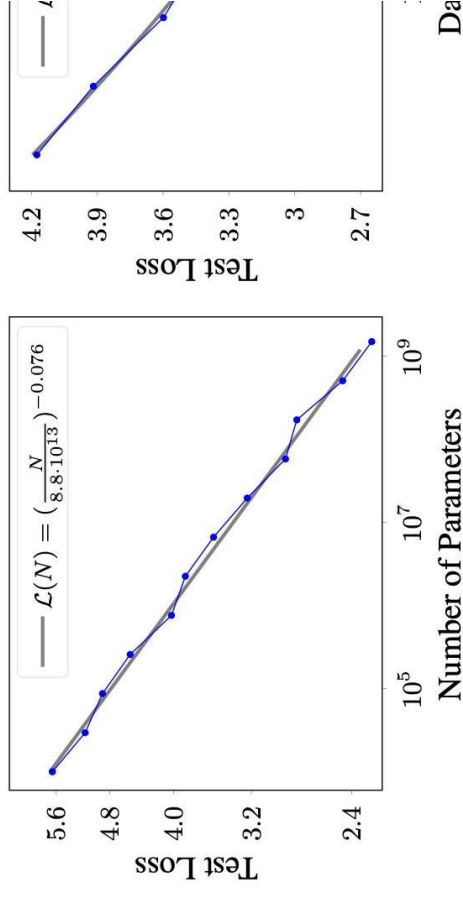
Why (1/2): self supervised/ scaling law

1 - Self-Supervised Learning pre-training



Lample & Conneau 2019

2 - Scaling Law



Compared to supervised method (first slides) there is no need to label the data with classical method

=> No physics to assum

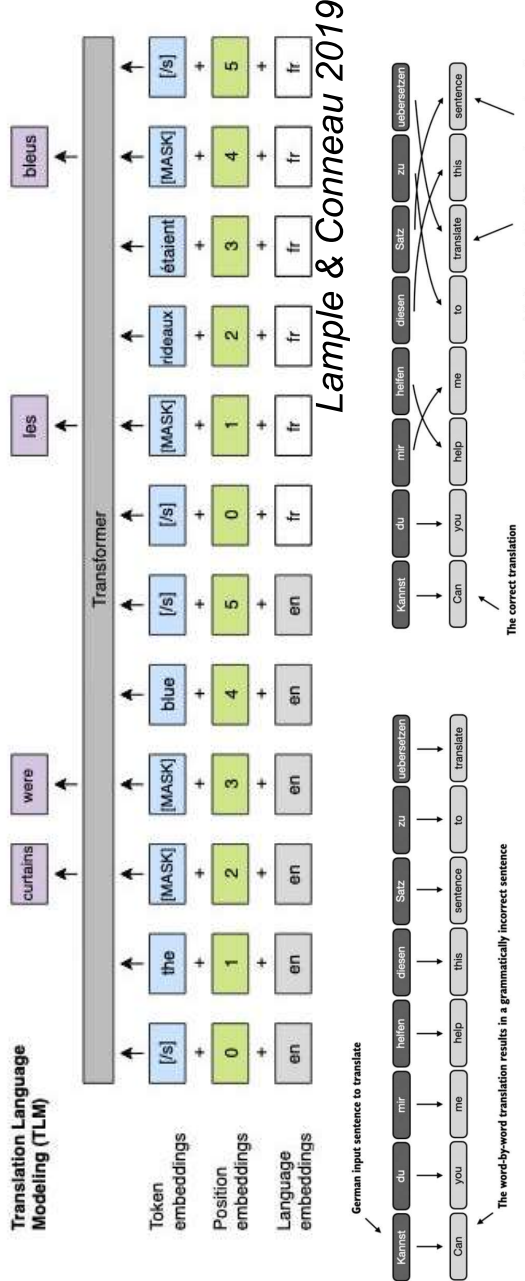
=> No need to first create label, you can use a larger dataset, important when label are hard to derive (ex: non LTE physics)

There is theroretically no limit or slowing in

=>Each night of WST will improve the

Why (2/2): Generalization / Versatility

3 - Cross-Domain Generalization



X Supervised training

✓ LLM
Sebastian Raschka

Due to self-attention mechanism and tokenisation, LLM learn universal patterns, not memorizing sequences and are permutation-invariant reasoner (e.g., "redshift = 0.5" and "z = 0.5")

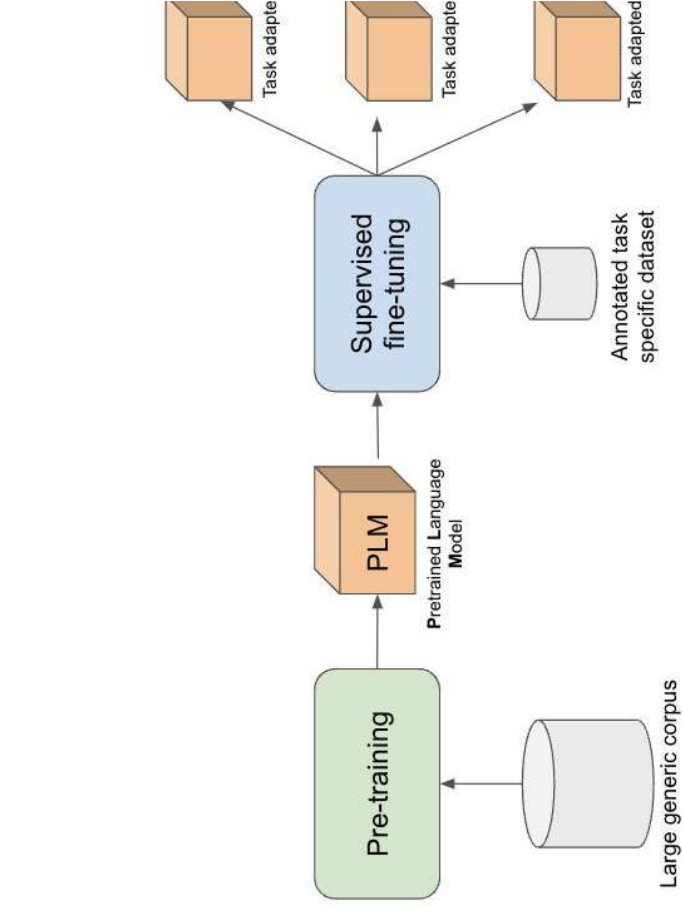
=>LLM can be trained on diverse observational data spanning cosmic scales (galaxy, stars...)
=>LLM can be trained on diverse instruments (eg: WST and MOONS)

AI & cINN

Applied to stars

LLM

4- Versatility



After pre-training you can finetune for very di with very few labels (100-1000)

Eg: create WST-like spectra, classify galaxy,

=> Foundation models arrive pre-trained a deploy, enabling diverse teams to implement AI capabilities with minimal technical ove

Foundation Model in Astro

‘Chat GPT’ like



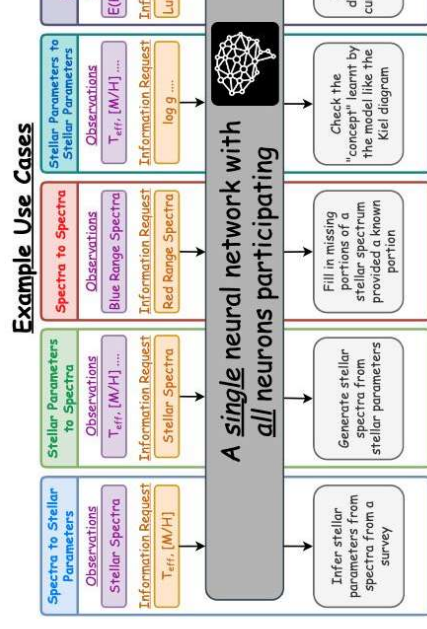
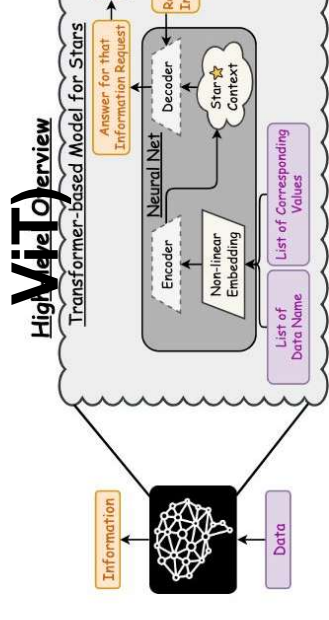
Nguyen et al.(2023)

AstroLlama



Cyuca & Ting 2023,

Science data processing



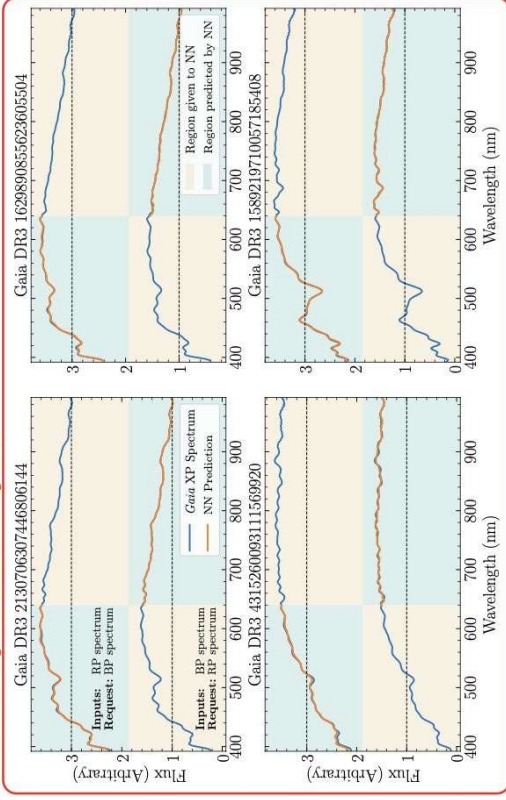
Credit: Brain icon by Imagisthanel and Astronomer icon by M...

Lung & Boy
Foundation Model
Stars with a
Transformer-based



LLM for Astro

Task: Stellar Spectra to Stellar Spectra



Dataset:

110 Gaia XP coefficients,

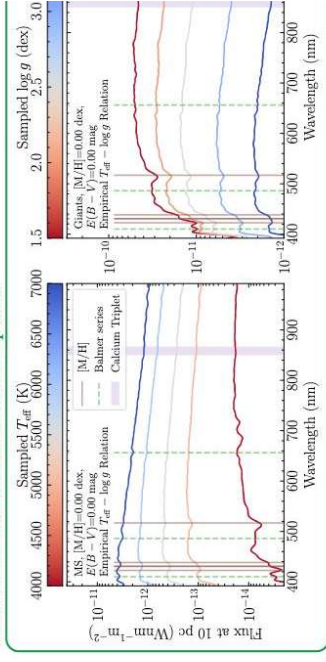
$G_{BP} - G_{RP}$, $J - H$, $J - K$,

T_{eff} , $\log g$, $[M/H]$,

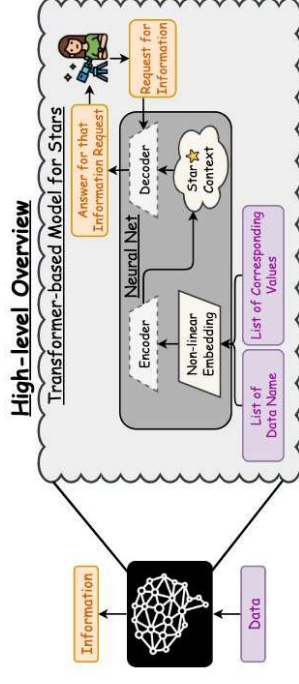
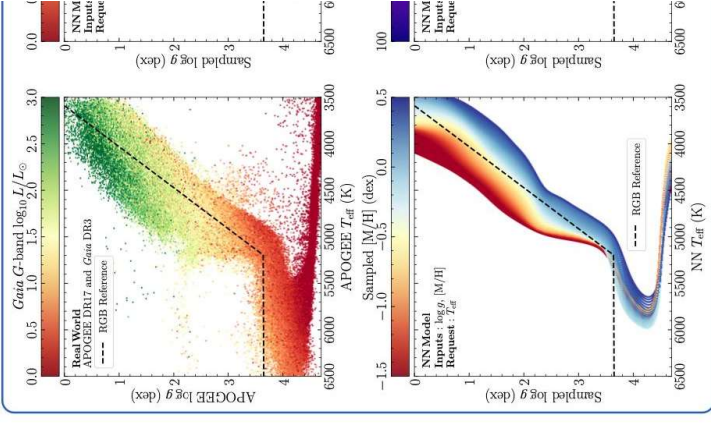
logarithmic $E(B-V)_c$,

and Gaia G -band pseudo-luminosity

Task: Stellar Parameters to Stellar Spectra

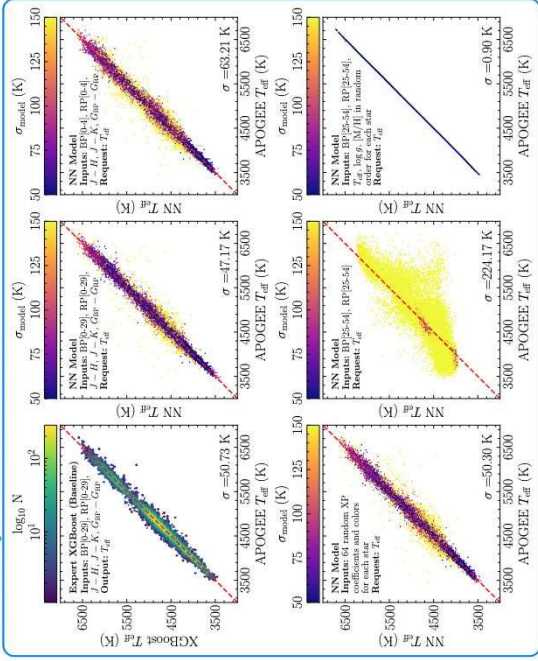


Task: Stellar Parameters to Stellar Parameters

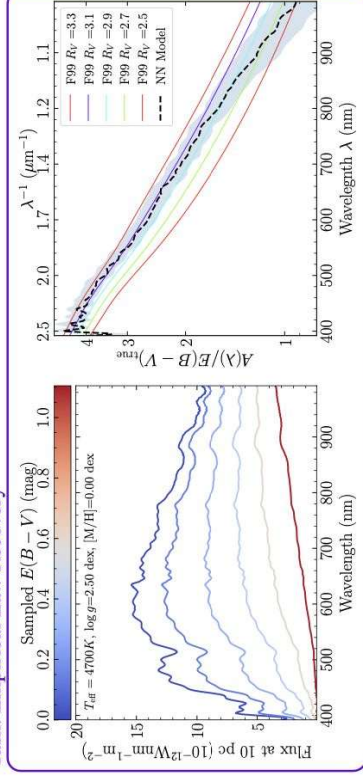


Lung & Bovy (2023)

Task: Stellar Spectra to Stellar Parameters



Task: Empirical Law Recovery



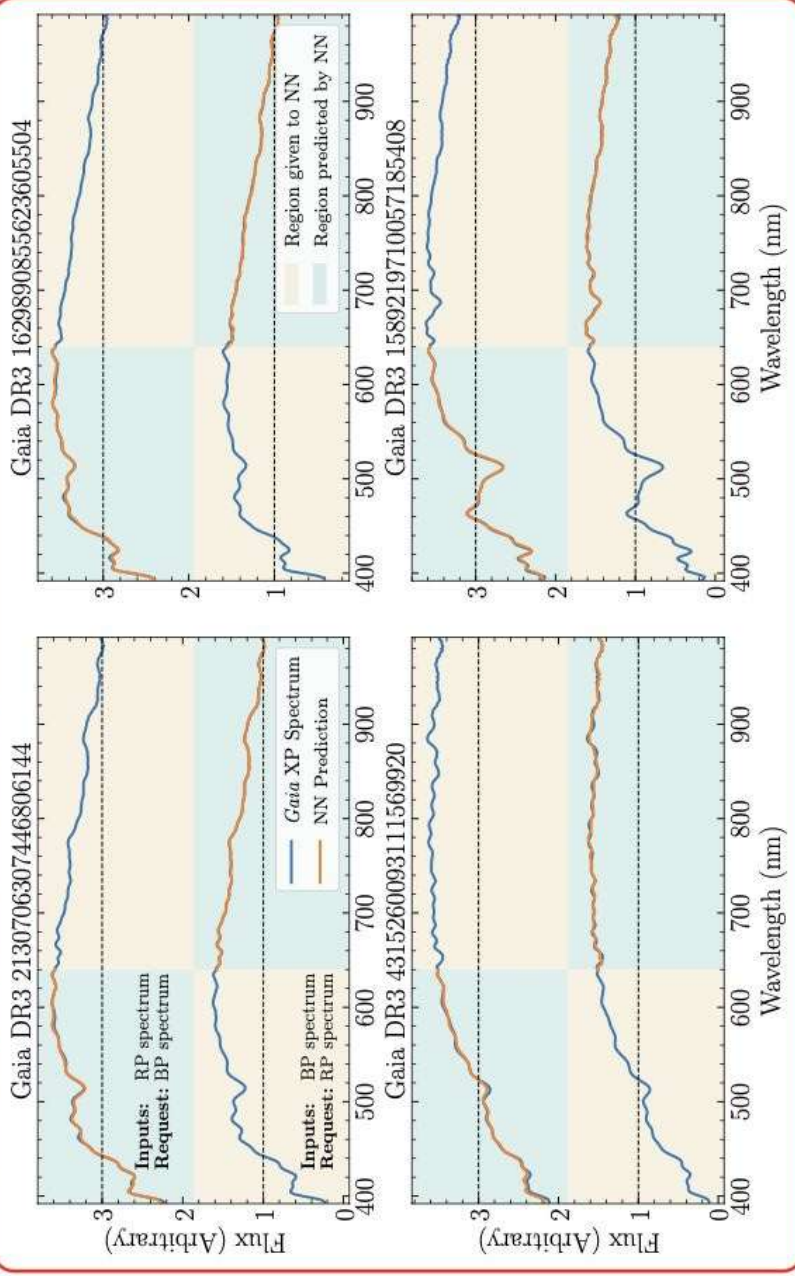
AI & cINN

Applied to stars

LLM

LLM for Astro

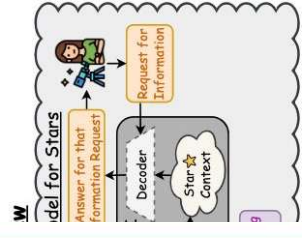
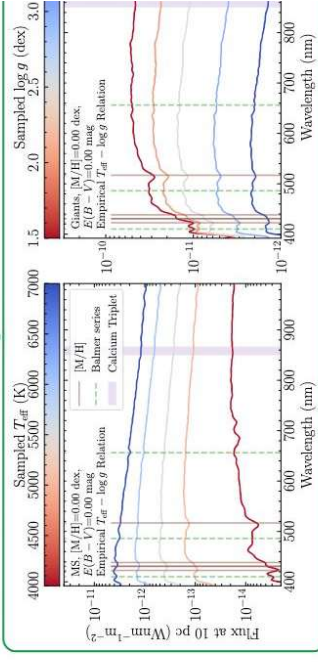
Task: Stellar Spectra to Stellar Spectra



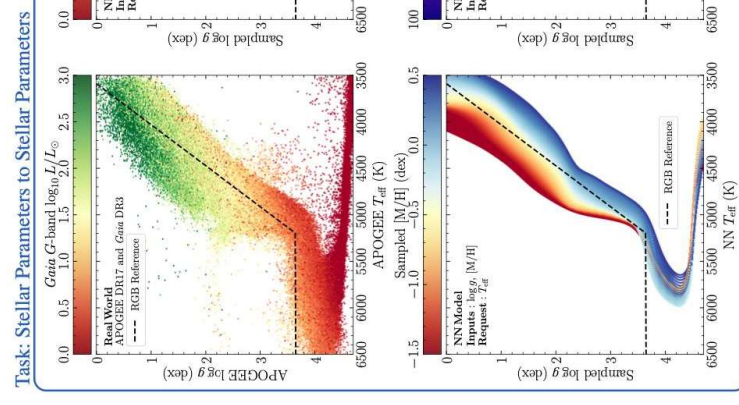
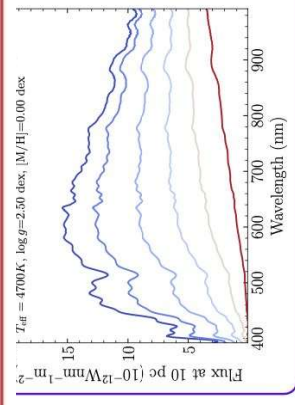
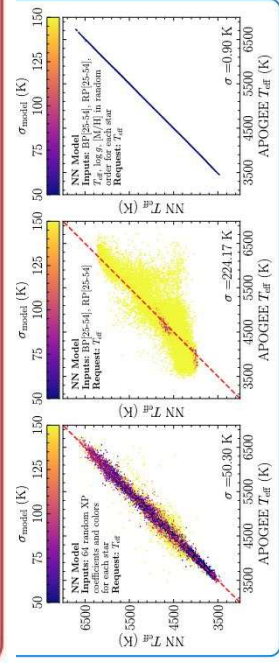
nts,

luminosity

Task: Stellar Parameters to Stellar Spectra



Bovy (2023)



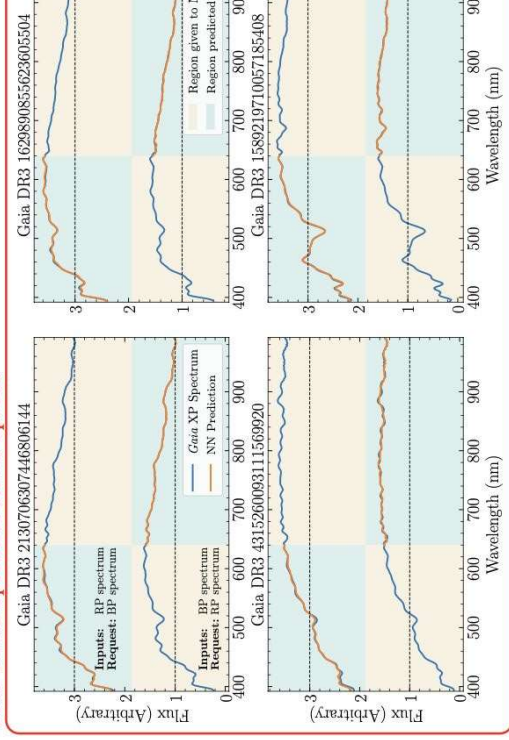
AI & cINN

Applied to stars

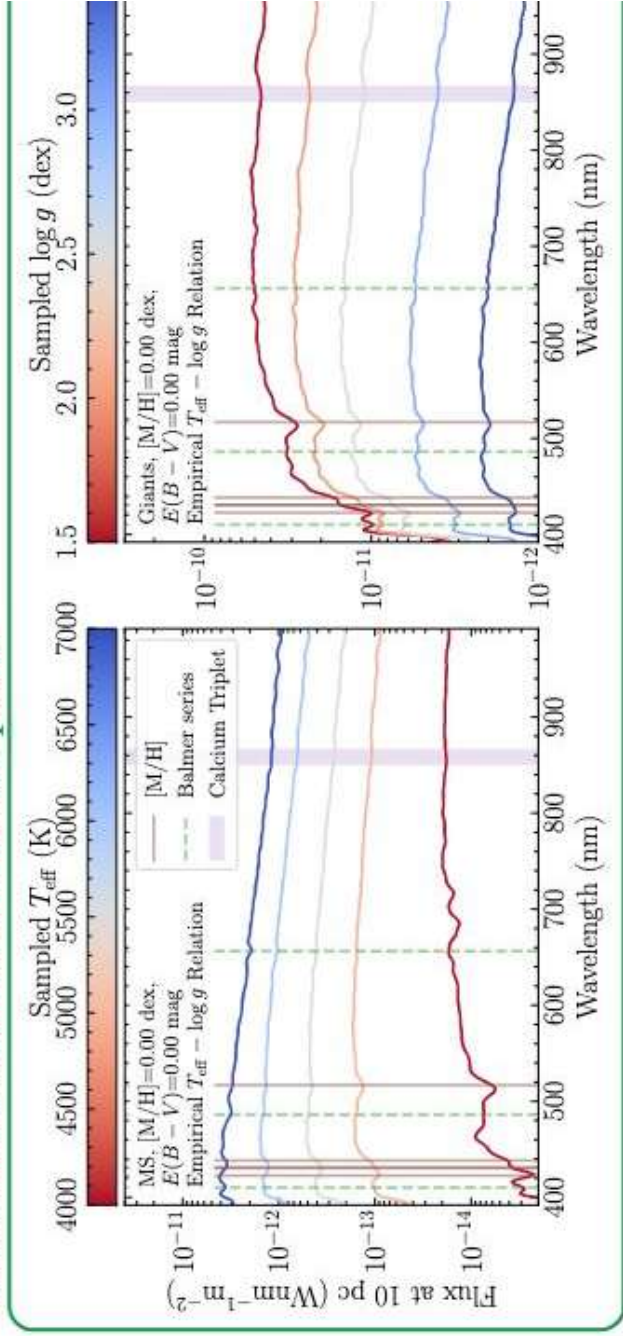
LLM

LLM for Astro

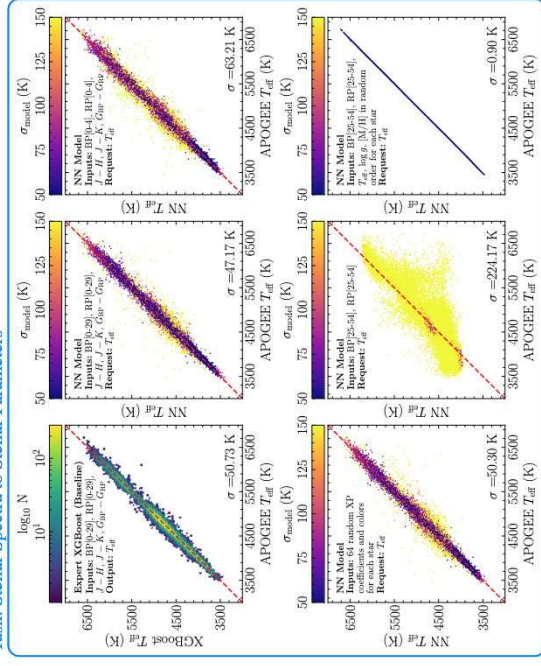
Task: Stellar Spectra to Stellar Spectra



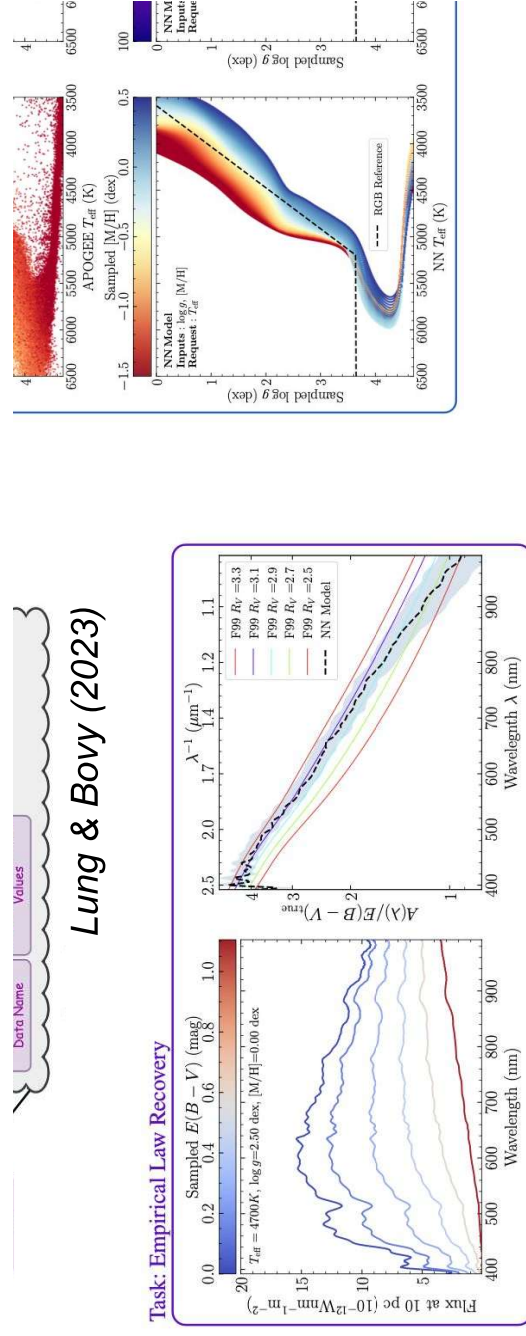
Task: Stellar Parameters to Stellar Spectra



Task: Stellar Spectra to Stellar Parameters

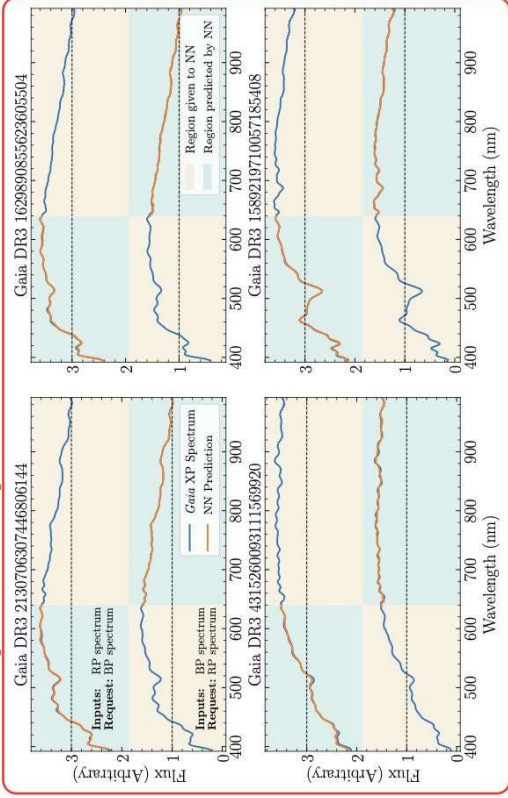


Task: Empirical Law Recovery



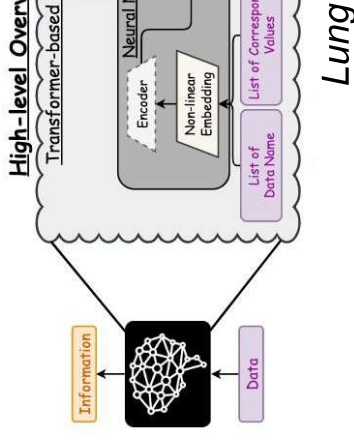
LLM for Astro

Task: Stellar Spectra to Stellar Spectra



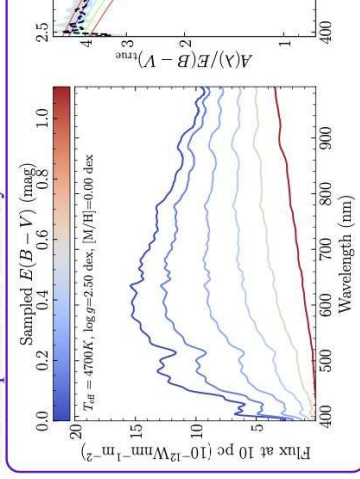
Dataset:

110 Gaia XP coefficients
 $G_{BP} - G_{RP}$, $J - H$, $J - K$
 T_{eff} , $\log g$, $[M/H]$,
 logarithmic $E(B-V)_c$,
 and Gaia G -band pseudo-color



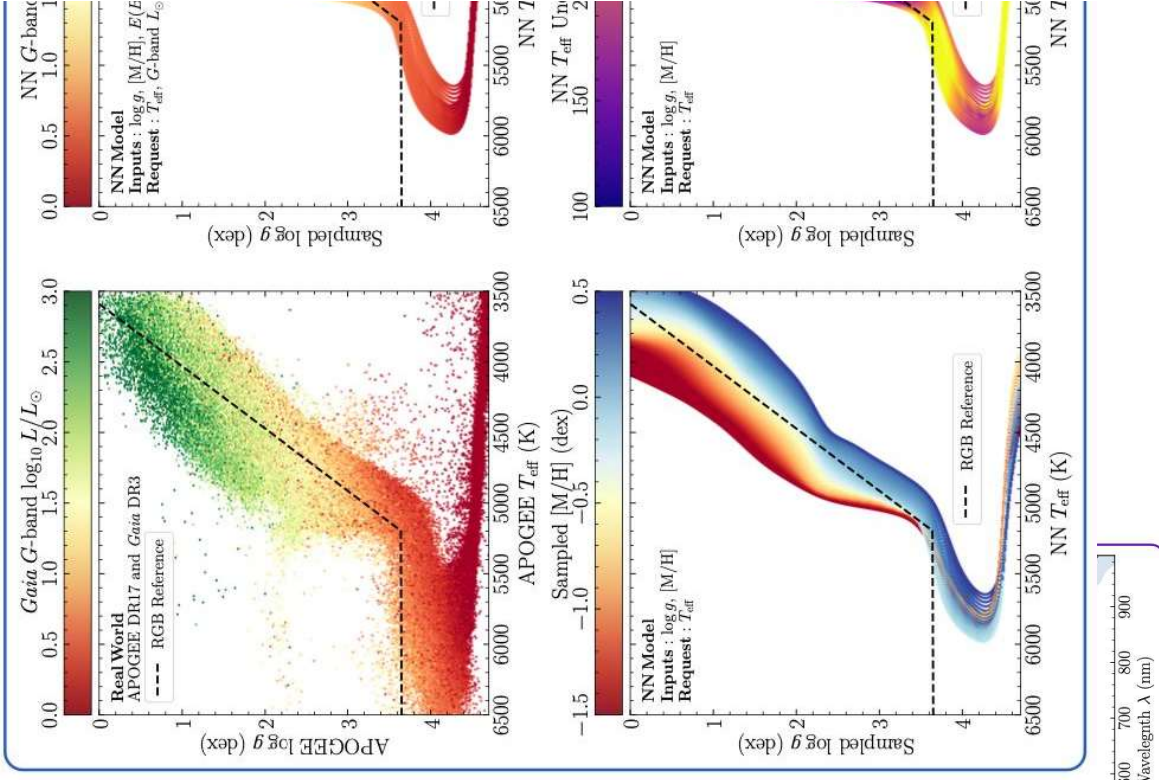
Lung

Task: Empirical Law Recovery



Task: Stellar Parameters to Stellar Spectra

Task: Stellar Parameters to Stellar Parameters



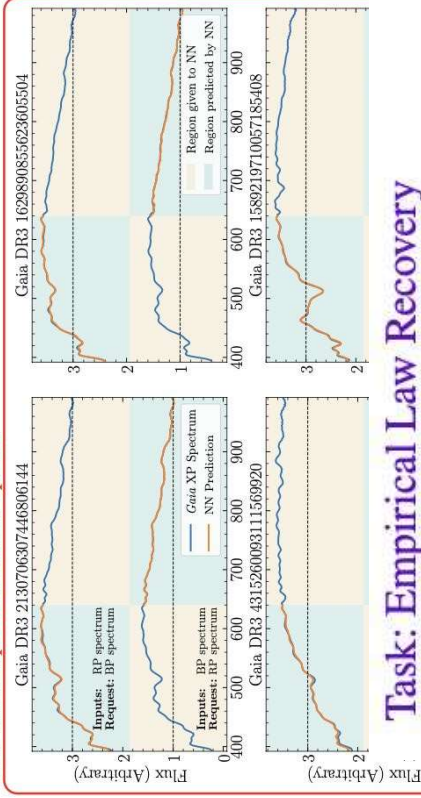
AI & cINN

Applied to stars

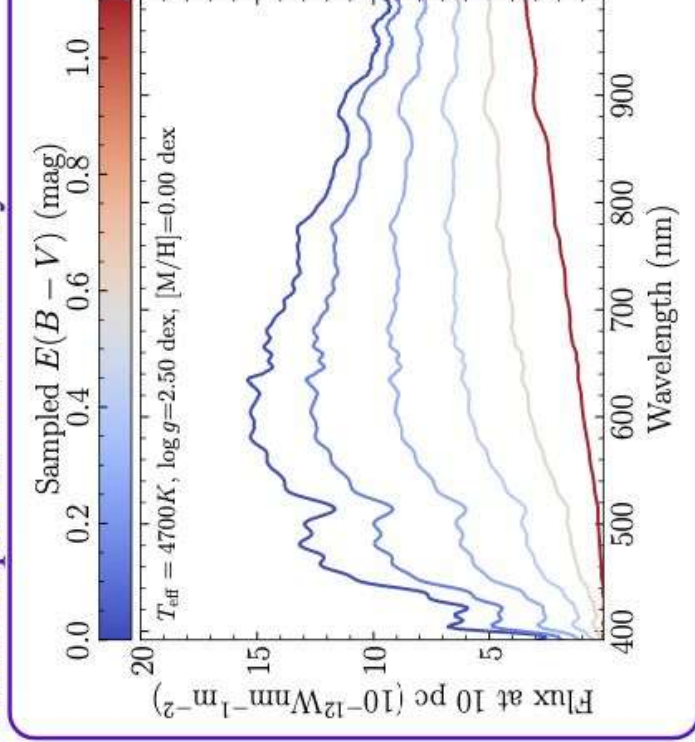
LLM

LLM for Astro

Task: Stellar Spectra to Stellar Spectra



Task: Empirical Law Recovery



Dataset:

Spectrum: 110 Gaia XP coefficients,

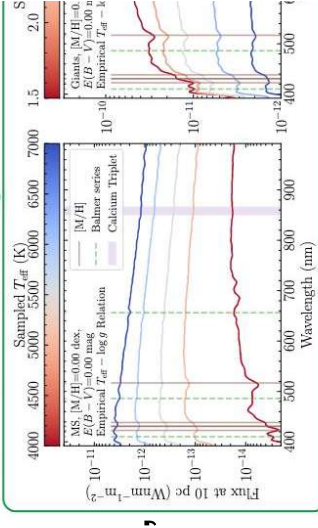
$G_{BP} - G_{RP}$, $J - H$, $J - K$,

T_{eff} , $\log g$, $[M/H]$,

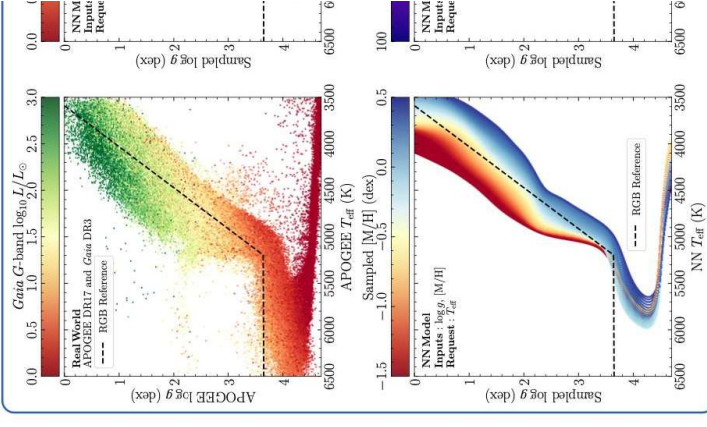
Logarithmic $E(B-V)_c$,

GaiaG-band pseudo-luminosity

Task: Stellar Parameters to Stellar Spectra

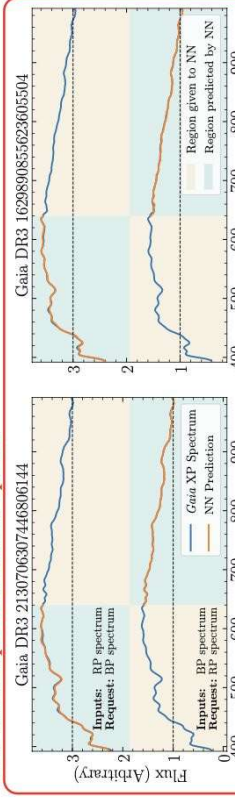


Task: Stellar Parameters to Stellar Parameters

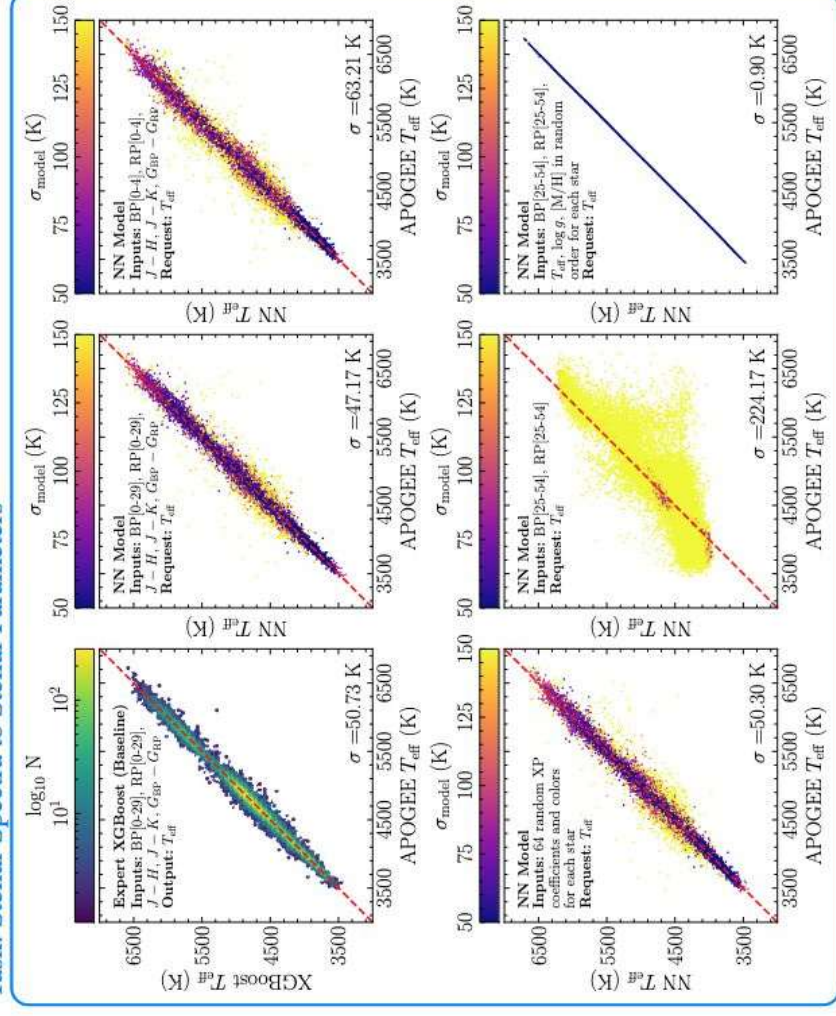


LLM for Astro

Task: Stellar Spectra to Stellar Spectra



Task: Stellar Spectra to Stellar Parameters



Dataset:

110 Gaia XP coefficients,

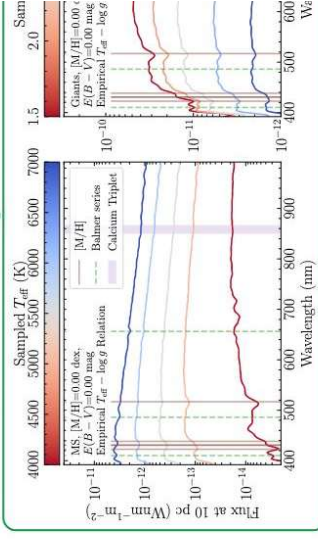
$G_{BP} - G_{RP}$, $J - H$, $J - K$,

T_{eff} , $\log g$, $[M/H]$,

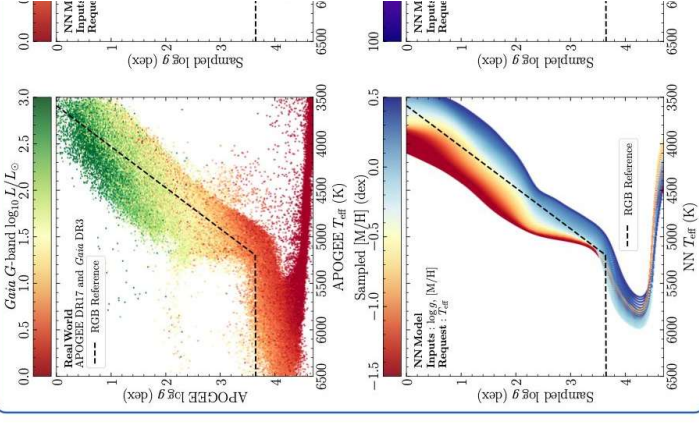
\ln arithmic $E(B-V)_c$,

$\tilde{\lambda}$ -band pseudo-luminosity

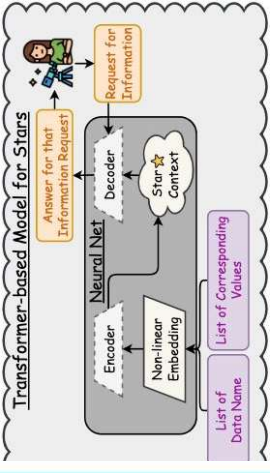
Task: Stellar Parameters to Stellar Spectra



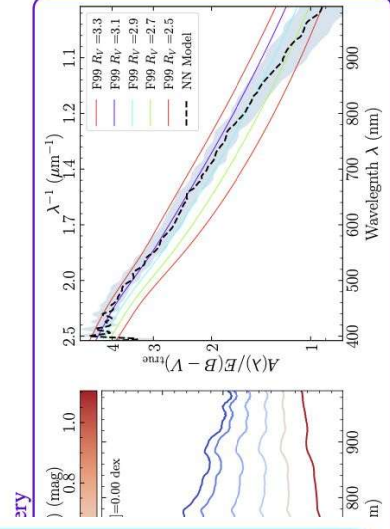
Task: Stellar Parameters to Stellar Parameters



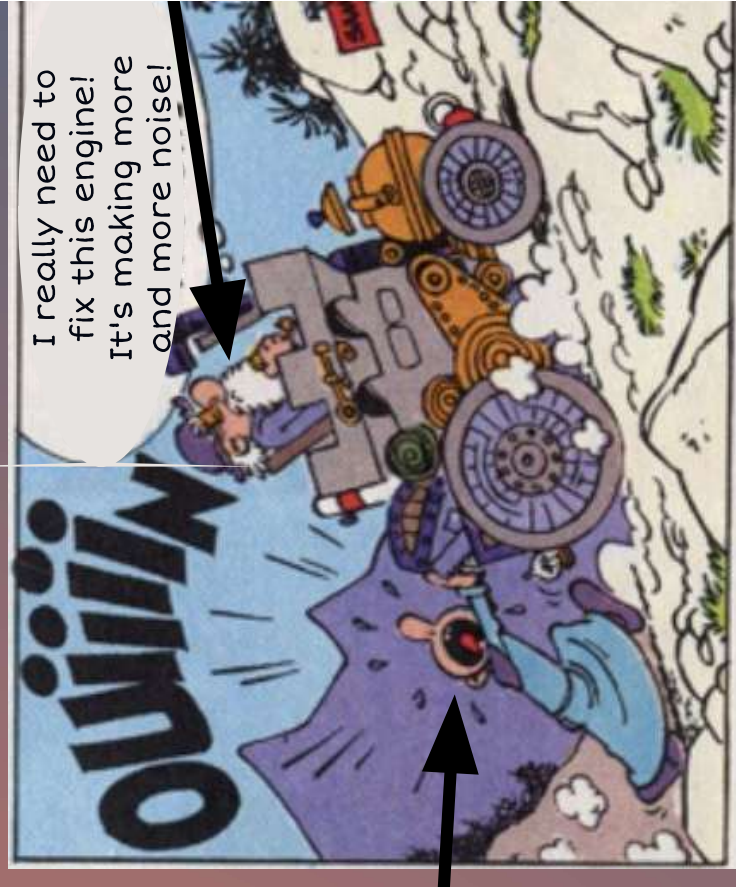
High-level Overview



Lung & Bovy (2023)



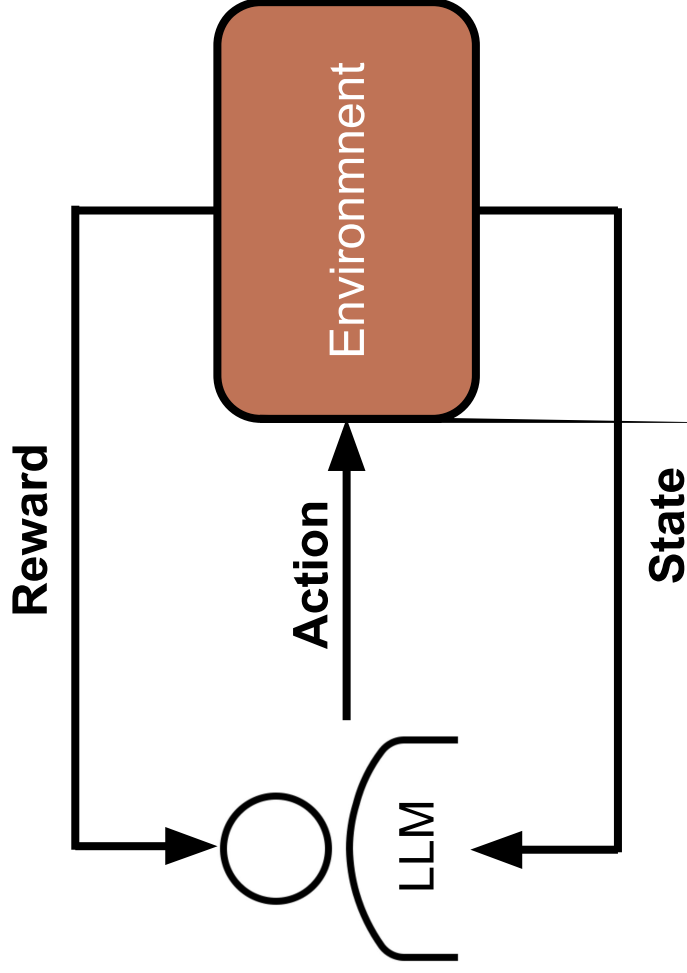
AI agent



AI Agent

Leonard, un genie en balade

What is an AI Agent



An artificial intelligence (AI) agent is a program that can interact with its environment, collect data, and use the data to perform self-determined tasks to meet predetermined goals.

Humans set goals, but an AI agent independently chooses the best actions it needs to perform to achieve those goals.



Can I have a cup of coffee please?

Sure! let me prepare it

Planning:

The user wants me to prepare coffee. I need to check if the coffee machine is connected and determine which coffee the user wants.

Action:

{Env: coffee mach}: Is the coffee machine connected? is the tank water empty?

{Env: clock}: what times is it?

State:

- The coffee machine is connected
- The tank water is full

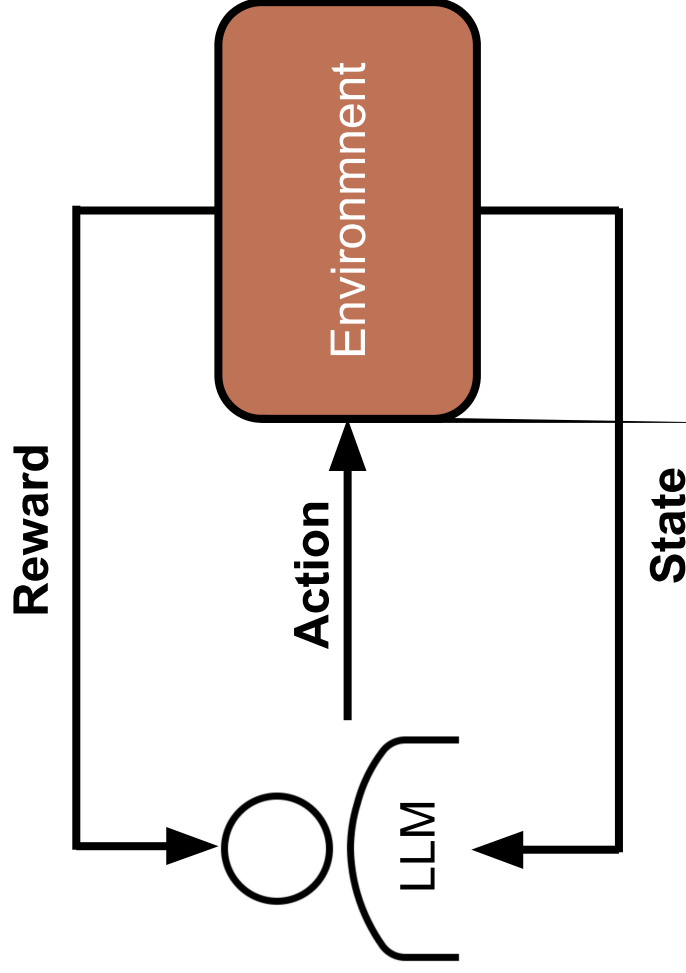
State:

- It is 15h30

Planning:

The coffee machine is connected and ready to prepare a coffee. It is 15h30, hence it is too late for a cappuccino. Let's go with a simple espresso

What is an AI Agent

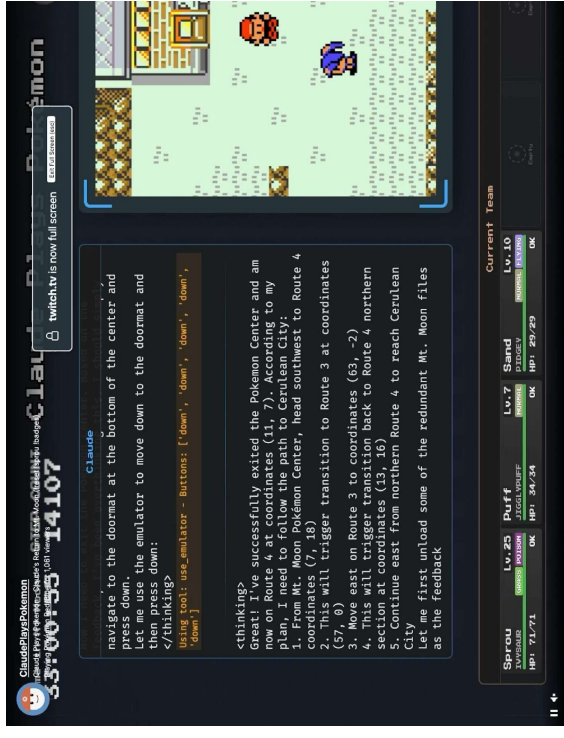
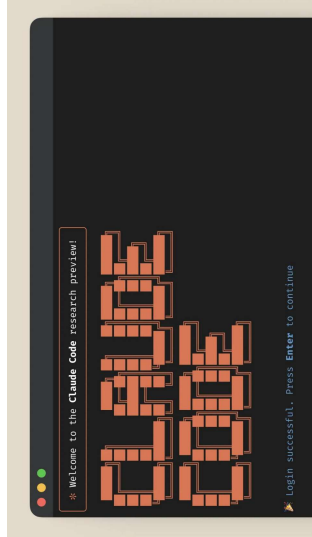
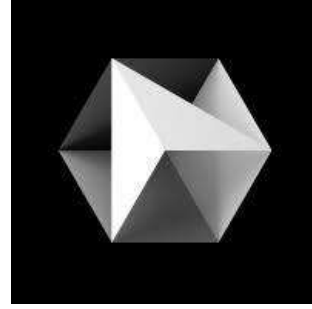


An artificial intelligence (AI) agent is a program that can interact with its environment, collect data, and use the data to perform self-determined tasks to meet predetermined goals.

Humans set goals, but an AI agent independently chooses the best actions it needs to perform to achieve those goals.

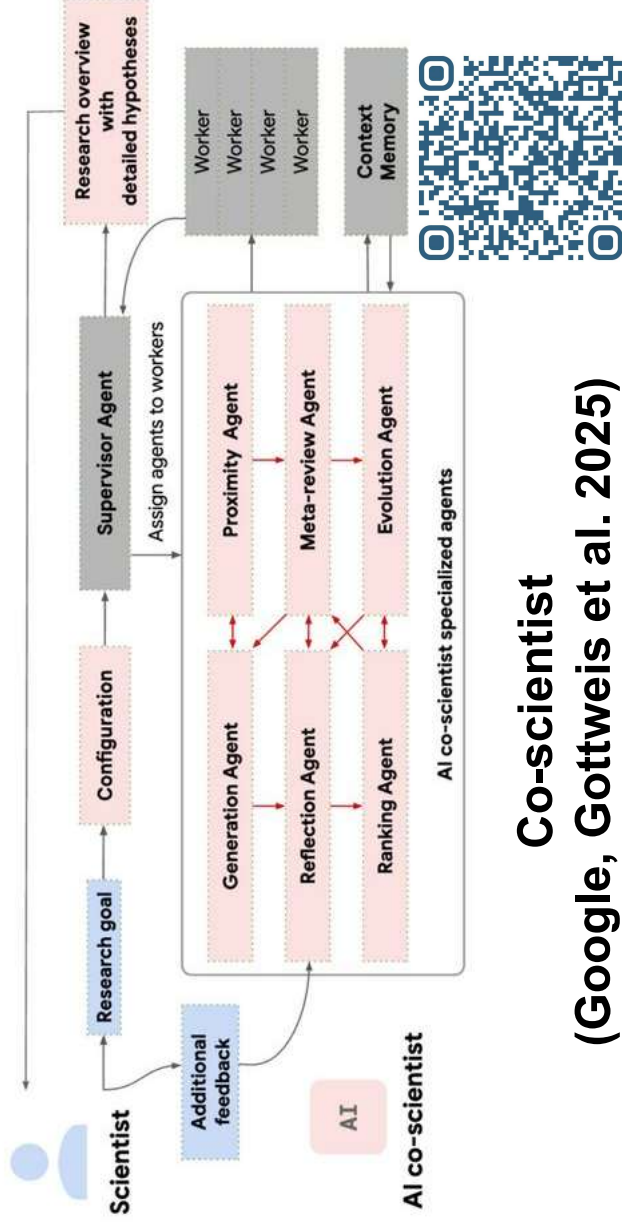


Example of agents



Coding

Playing Pokemon



Co-scientist
(Google, Gottweis et al. 2025)

AI & cINN

Applied to stars

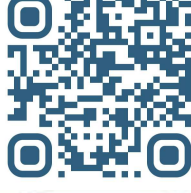
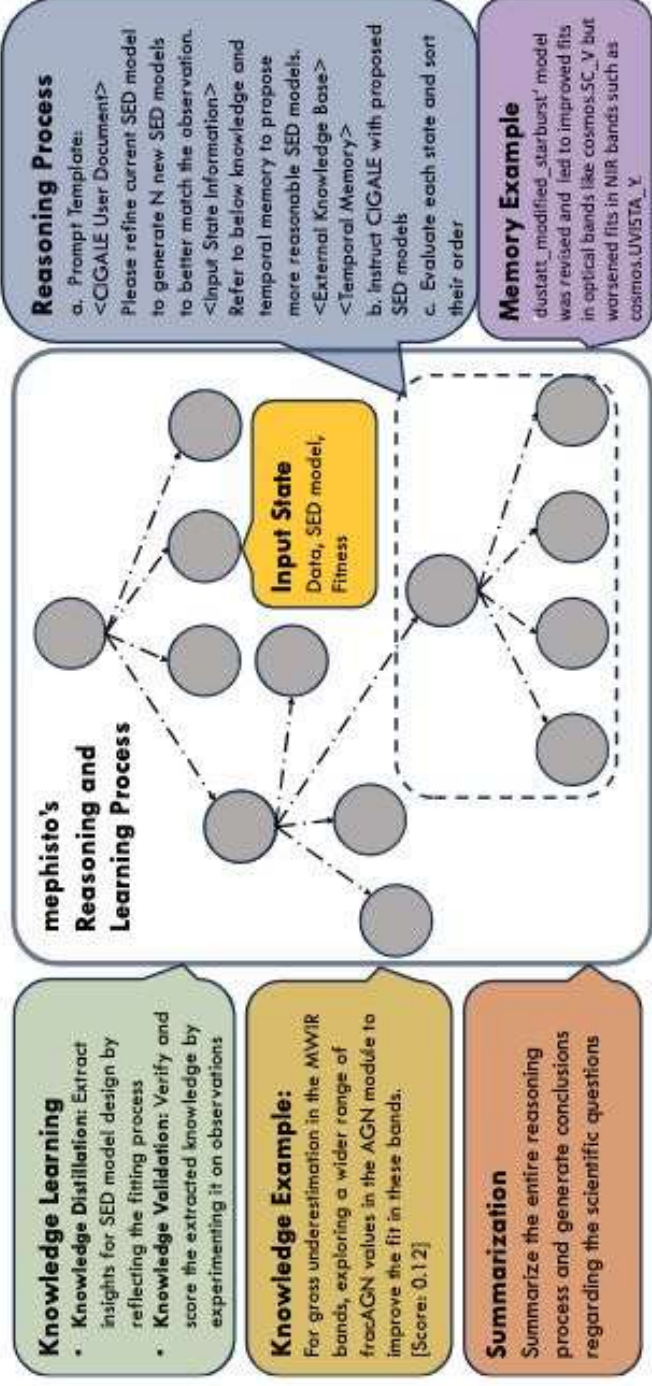
LLM

AI Agent

Data Science Gemin

Agent AI Astro

Mephisto, an AI agent able to use CIGALE code to fit SED



Sun et al. 2024.

In 2040, AI agent will be use both for the processing of data and for the planning of observation (target decision, mode choice, opportunity decision...)

<CIGALE Documentation>

<Introduction to User Input>

Your task is to analyze user input and adhere to below rules: provided CIGALE models to improve the fit quality of different the CIGALE module which should affect the fitting quality the module, e.g., use another choice, adjust the parameter grid, and

- **Choice Selection:** For each module in the model, select the best choice based on the need to optimize observational data.
- **Parameter Grid Specification:** For the selected choice, provide a parameter grid that the model will use to fit the data.
 - For discrete parameters, select grid values from the provided range
 - For continuous parameters, derive grid values with a comprehensive exploration of parameter space
- **Mandatory Modules:** The model configuration must include the following modules:
 - sfb (Star Formation History): This module is essential for modeling a galaxy forms stars over time.
 - ssp (Simple Stellar Population): This module is used to model the collective properties of stars in a galaxy that form from a single population of stars with the same metallicity.
- For 'imf' parameter in 'bc03' and 'm2005', it should be '1' are acceptable, 'imf': [1] are forbidden
- For 'disk_type' parameter in 'fritz2006' and 'skirtor', it should be '0' are acceptable, 'disk_type': [0]

##CIGALE SED Knowledge##

To enhance the fit quality for different photometric bands in the following additional information when designing your CIGALE model:

##Temporary Memory##

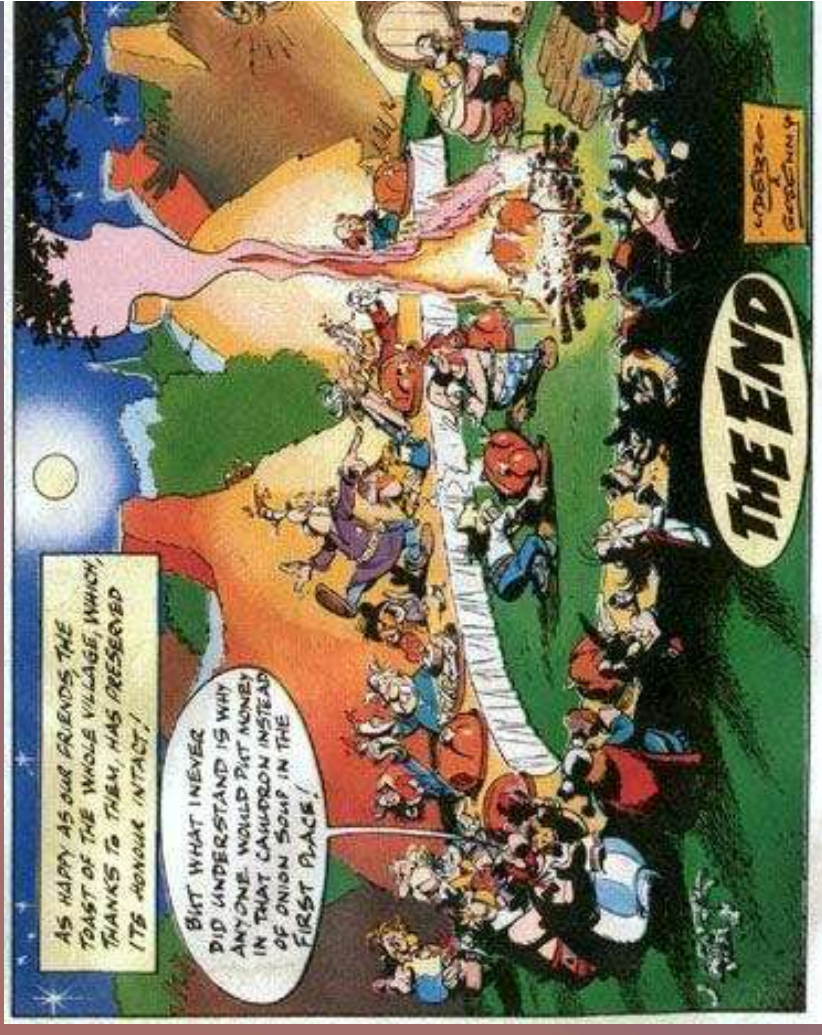
Please take into account the following temporary memories to improve your fit quality in the updated CIGALE model:

##Expected Output Format##

Your output format should be structured as below list constructible CIGALE model modifications:

```
[
  {
    "thinking": "thinking for CIGALE model",
    "module": "module name for CIGALE model",
    "name": "module choice name for CIGALE",
    "parameters": [
      "parameter 1": ["parameter 1 grid h", "parameter 1 grid b", ...],
      "parameter 2": ["parameter 2 grid h", "parameter 2 grid b", ...],
      "parameter n": ["parameter n grid h", "parameter n grid b", ...]
    ]
  },
  ...
]
```

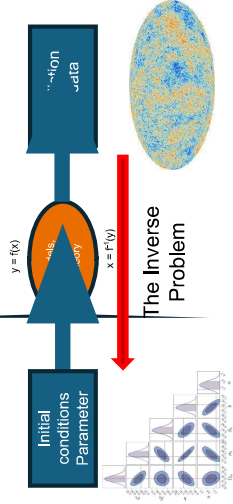
Conclusion



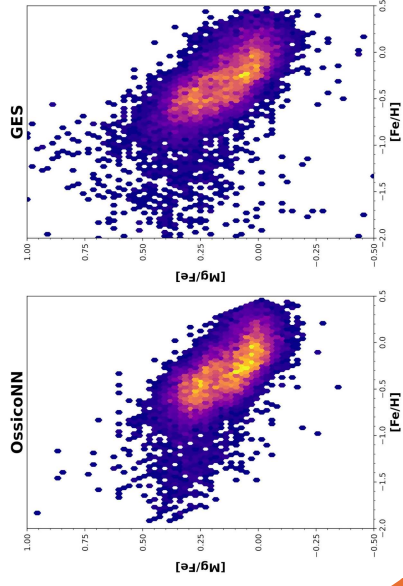
Asterix and Obelix and the cauldron

Conclusion

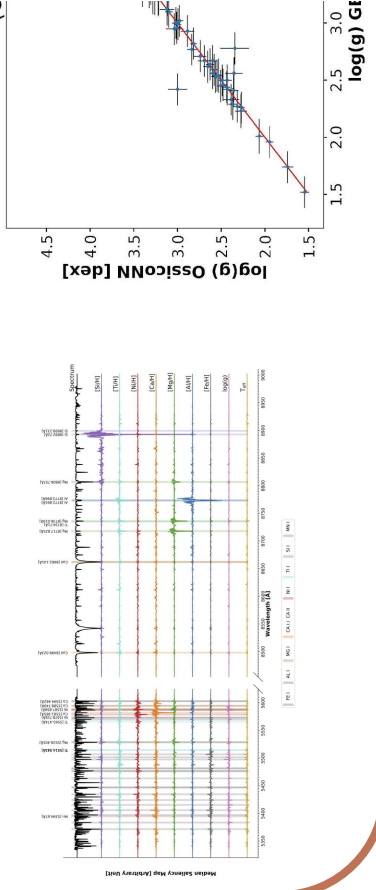
AI is important to tackle inverse problems (but not exclusively)



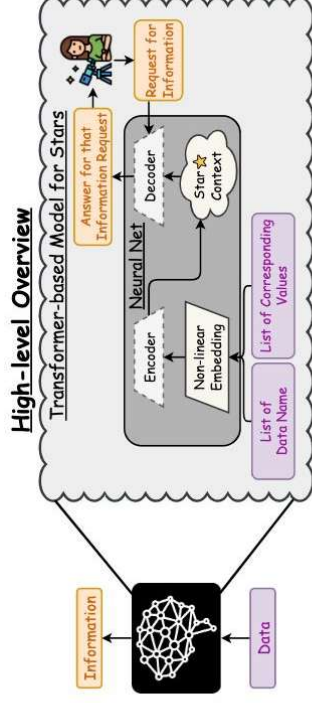
AI well reproduce Astrophysical relation



NOT A BLACK BOX: we get feature importance and comprehensive uncertainty



In 2040, LLM will be used both for the processing of data and communication of results



By 2040, astrophysicists will primarily be planners, relying on AI agents to extract information from observational tasks and data analysis



Agent



Thanks !!!!

