# Astronomicky velká Velká data
# Virtuální observatoř
# a AI v astronomii

## Petr Škoda

Astronomický ústav AVČR

Setkání Kosmologické sekce ČAS
Matematický ústav AVČR, Praha, 11.3.2026

# Credits

- The presentation is based on many different sources – mainly the on-line published slides from IVOA meetings, slides from Astroinformatics workshops, ADASS conferences, COST Action TD1403 or pictures found on Internet.

# Outline of the Talk

- Data Avalanche in astronomy

- Virtual Observatory

- Astroinformatics (AI)

    - Visualizations

    - Transfer of technology
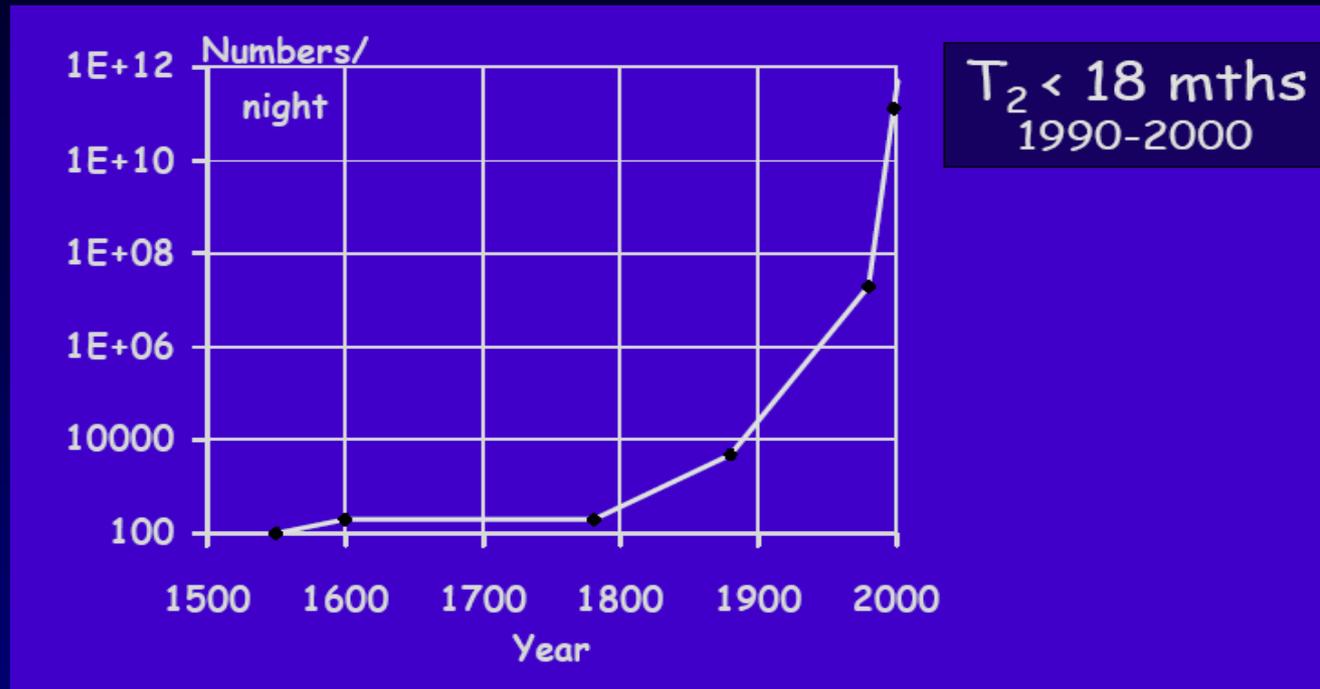
- Artificial Intelligence (AI)

- Future?

# *Astronomically*

# *Big Data*

# Data Avalanche

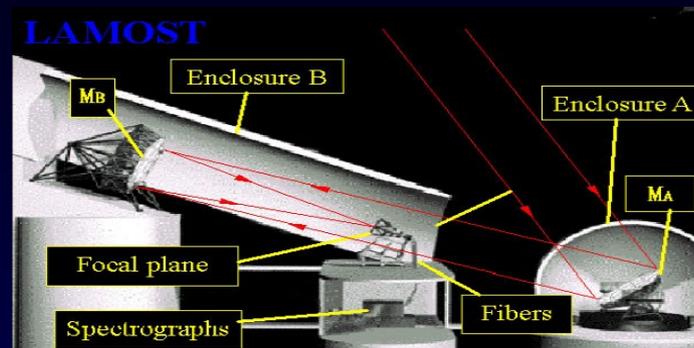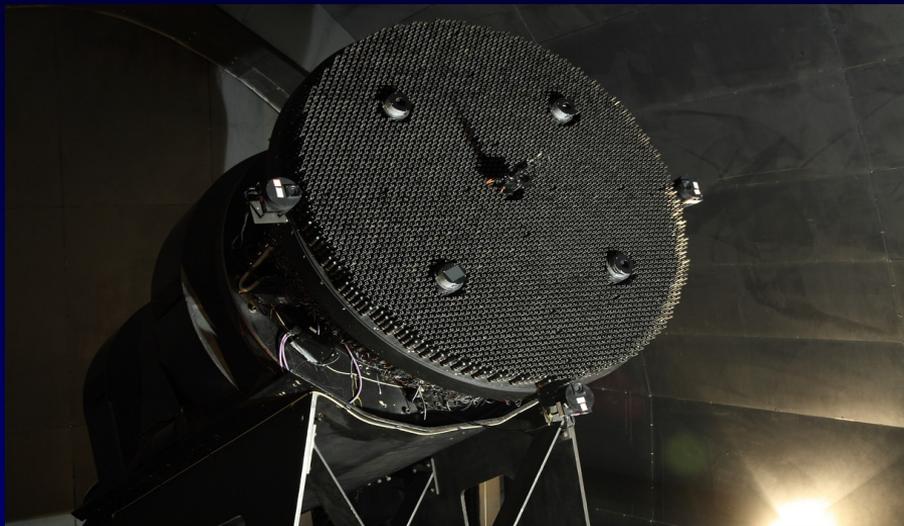Moore law for chips –doubling 1.5 year

Data in astronomy – doubling < 1 yr !

100 PB today,   100 TB/night

# LAMOST (Guoshoujing)

- Xinglong, China
- 4 m mirror (30 deg meridian)
- 4000 fibres

# LAMOST Spectral Surveys

DR1 (end 2013)        **2 204 860**  spectra
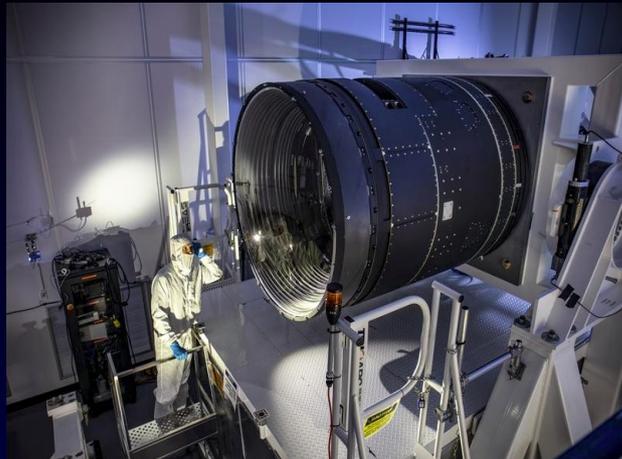                              1 085 404 stars classified by pipeline
DR2 (beg 2015)      **4 132 782** spectra
                              3 779 674 stars
                              307 000 unknown!
DR12  (Mar 2025)      12 602 390 low res
                            + 15 475 985 mid res

Each fibre – 2 motors
double arm 33mm circle

Fibre collects light from
3.3 arcsec circle on sky

# LSST – Vera C. Rubin Observatory



189 CCD  4kx4k, 10um
3.2 Gpix every 15 sec
3.5 deg FOV (64cm)
15 TB/day=6 PB/yr RAW
15 PB catalogue (D11)
detection of changes 60s!
7 million allerts/night ! Tot 20B
38 billion objects x 800
32 tril. meas. -5 PB table

- **Virgo cluster in Rubin**
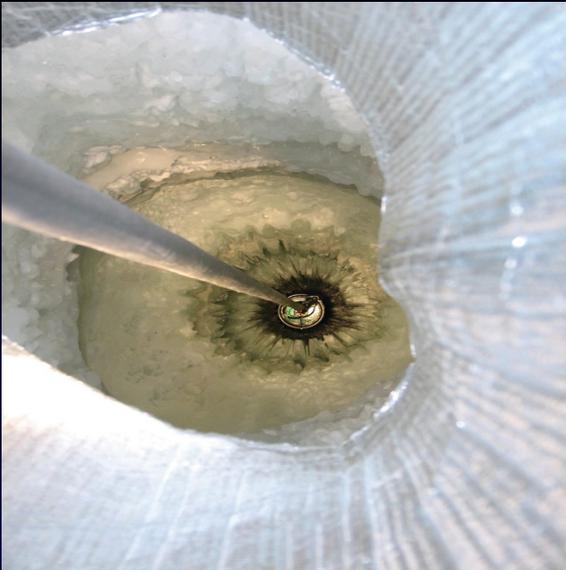
  **15 sq deg**

  **1185 exposures / 7 nights**

Credit: NSF–DOE Vera C. Rubin Observatory/NOIRLab/SLAC/AURA

# Euclid



- 1.2m Korsch tel.
- Dark matter
- 10B sources
- 1B weak lensing
- 600 MP camera 7filt
- Spectra VIS/NIR
- DR1 on June 2026
- 100 PB processed
- 26PB/year



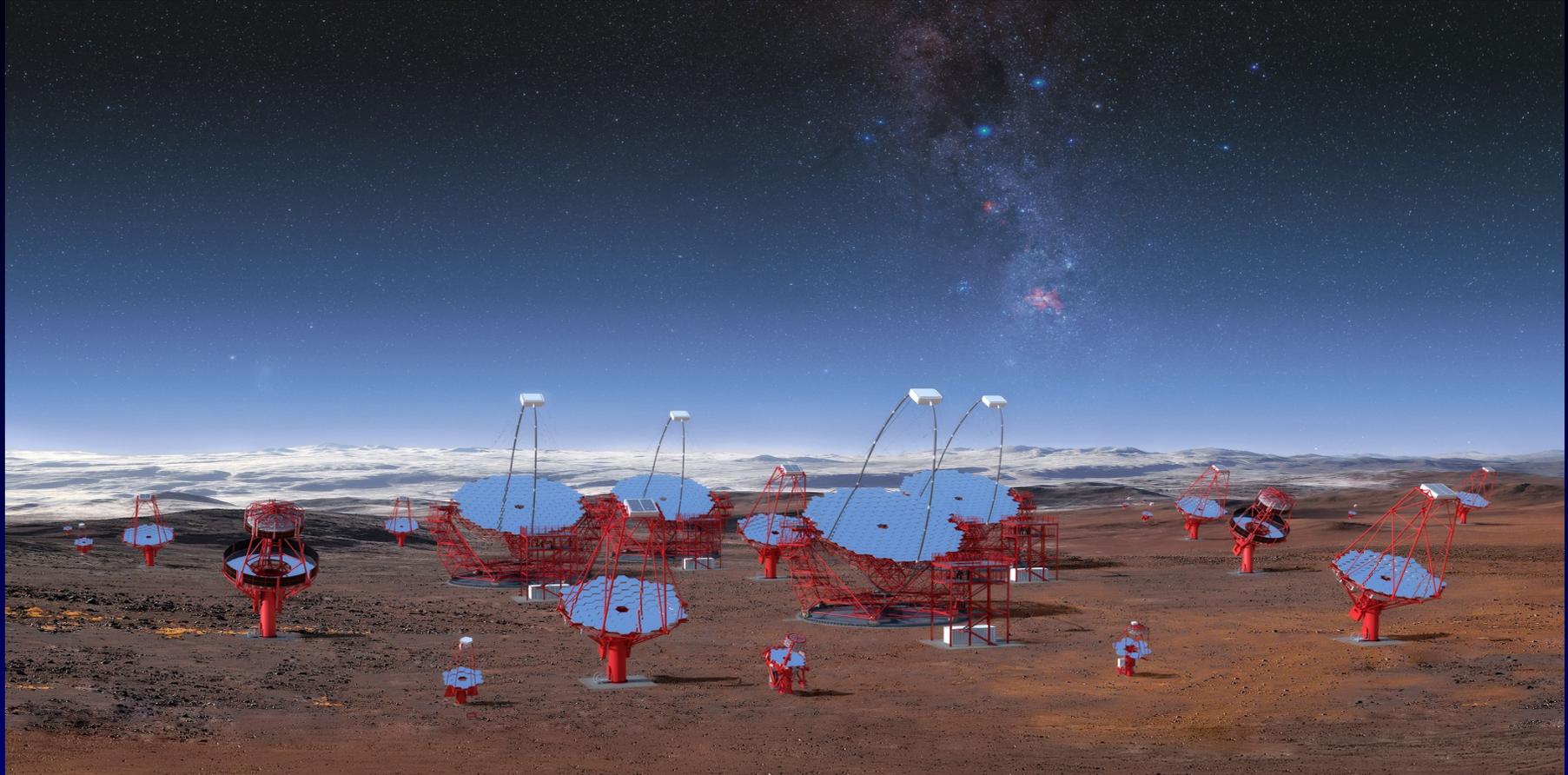Euclid Deep Field South J041110.98-481939.3, Credits: ESA/Euclid/Euclid Consortium/NASA,

# IceCube Neutrino Lab



South Pole
Amundsen-Scott station

# IceCube Neutrino Lab

# Cherenkov Telescope Array
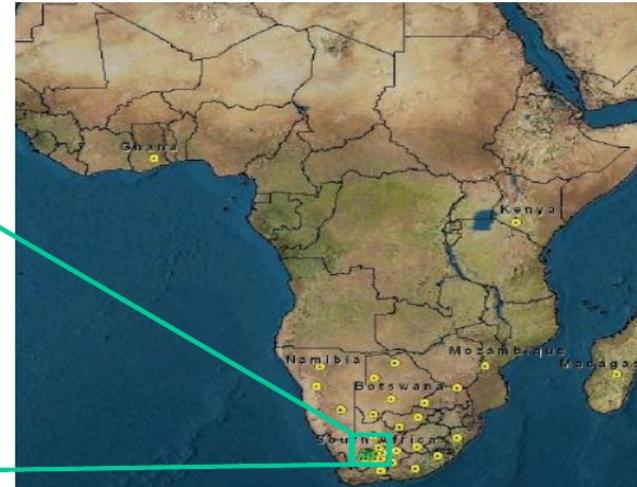
# Gravitation Wave Detection Network
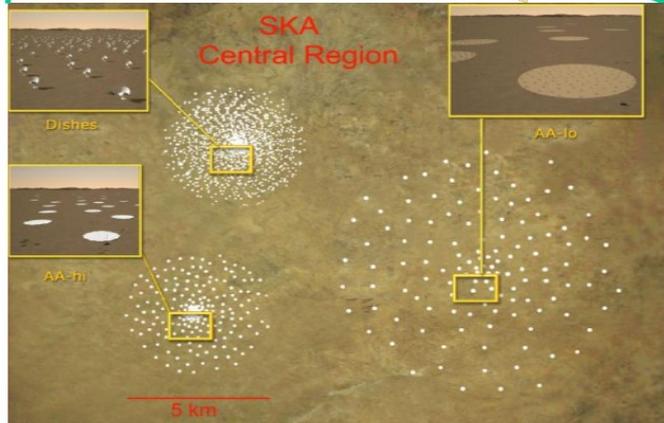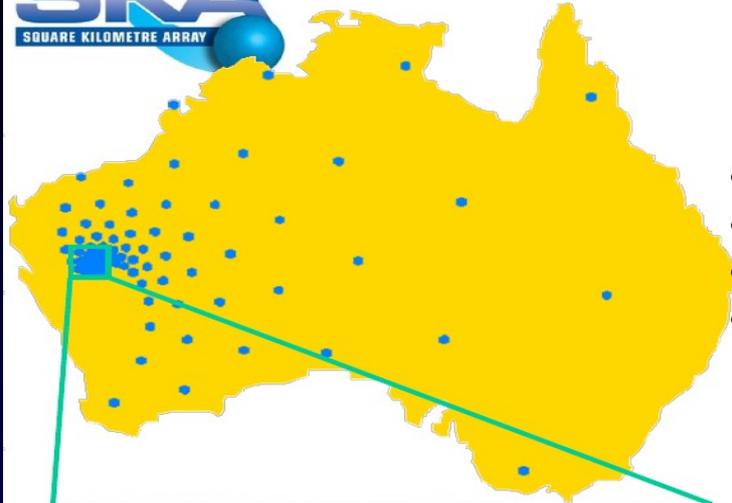
# LOFAR network



|  | LOFAR | SKA |
|---|---|---|
| Raw Telescope | 112 PB/yr | 60 EB/yr |
| Archive Rate | 6 PB/yr | 100 PB/yr |

# SKA



also a Continental sized Radio Telescope

- Need a radio-quiet site
- Very low population density
- Large amount of space
- Possible sites (decision 2012)
  - Western Australia
  - Karoo Desert RSA

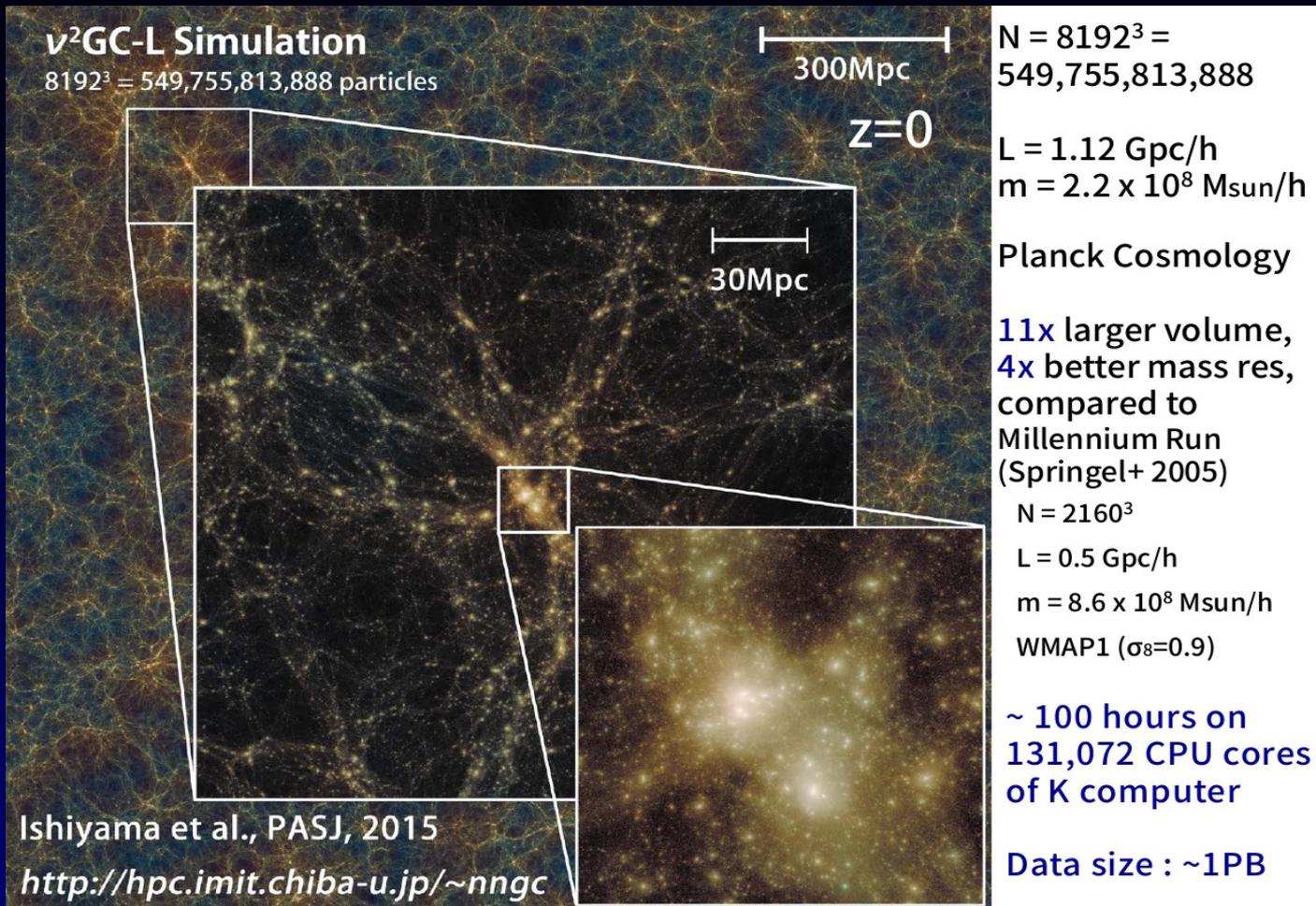# SKA



Dishes

# SKA



Phased Aperture array

# SKA Archive Volumes

- ~0.5 – 10 PB/day of image data
- Source count ~$10^6$ sources per square degree
- ~$10^{10}$ sources in the accessible SKA sky, $10^4$ numbers/record
- ~1 PB for the catalogued data

**100 Pbytes – 3 EBytes / year of fully processed data**

700PB/year   2x135 Pflops   20 Tbit/s transfer link

# Simulation of Universe



ν²GC-L Simulation
$8192^3 = 549{,}755{,}813{,}888$ particles

300Mpc

z=0

30Mpc

Ishiyama et al., PASJ, 2015
http://hpc.imit.chiba-u.jp/~nngc

$N = 8192^3 = 549{,}755{,}813{,}888$

$L = 1.12$ Gpc/h
$m = 2.2 \times 10^8$ Msun/h

Planck Cosmology

**11x** larger volume,
**4x** better mass res,
compared to
Millennium Run
(Springel+ 2005)
  $N = 2160^3$
  $L = 0.5$ Gpc/h
  $m = 8.6 \times 10^8$ Msun/h
  WMAP1 ($\sigma_8 = 0.9$)

~ 100 hours on
131,072 CPU cores
of K computer

Data size : ~1PB

# *Virtual Observatory*

# IVOA (established 2002)

# Virtual Observatory : Key Definitions

- *"The Virtual Observatory will be a system that allows astronomers to interrogate multiple data centers in a seamless and transparent way, which provides new powerful analysis and visualization tools within that system, and which gives data centers a standard framework for publishing and delivering services using their data".*

- Standardization of data and metadata, and of data exchange methods.

- Registry, listing available services and what can be done with them.

*R.J.Hanisch, P.J.Quinn, in "IVOA – Guidelines for participation"*

# Ecosystem of VO



**LEVEL 1 empty**

USERS

COMPUTERS

REC

InProgress

USER LAYER

Browser Based Apps — Desktop Apps — Script Based Apps

USING

| REGISTRY | Semantics | VO Query Languages | | Data Models | DATA ACCESS PROTOCOLS |
| --- | --- | --- | --- | --- | --- |
| | | VO CORE | | | |
| | | Formats | | | |

SHARING

Data and Metadata Collection

Storage — RESOURCE LAYER — Computation

20101004 IVOA Architecture

PROVIDERS

# Ecosystem of VO – level 2

# Spectra in SPLAT-VO - query

# Spectra in SPLAT-VO direct access

# Spectra in SPLAT-VO - DataLink

# TOPCAT

# Aladin

# Technology of VO

Unified data format– VOTable, UCD (Vizier)

Transparent transport  (VOunits)

VOregistry  (DNS like)  Google for data+WS

Protocols

    ConeSearch (searching in circle on sky)

    SIAP (Simple Image Access Protocol)

    SSAP(Simple Spectral Access Protocol)

    SLAP(Simple Line Access Protocol)  - VAMDC

    TAP (Table Access Protocol) – query e.g. whole SDSS

    VOEVENT (transients, robotic telescopes,Sun

    DATALINK (related data products, e.g. raw, mosaics..)

    SODA Server-side Operations for Data Acces

# Technology of  VO

ADQL (Astronomical Data Query Language)

    XMATCH, REGION  (2 catalogues – shifted)


Application interoperabilty  –  SAMP

    Allows develop applications as bricks

    sending VOTABLES  (catalogue-spectra-images)


Surveys visualization

    HIPS (Hiearchical Progressive Survey )  - allsky  zoom

    MOC  (Multi order coverages) time, space, spectral (FoV)

# Open Science – EOSC

EURO-VO DCA, ICE, CoSADIE, ASTERICS, ESCAPE - Astroparticles

# VO Science Portals



VO embedded in astronomy services

ESO Science Portal

WWT

Firefly
Caltech-IPAC

ESA Sky

Grav. waves

CDS reference data service

SVO Filter Profile service

# VO Science Portals

Stellarium + VirGo (ESO, unsupported)

ESASky
https://sky.esa.int/

ESO Archive Science Portal
https://archive.eso.org/scienceportal/home

IRSA IPAC archive (Firefly)
https://exoplanetarchive.ipac.caltech.edu/firefly/

WWT  (original MS, now AAS, web client)
http://worldwidetelescope.org/webclient

GoogleSky
https://www.google.com/sky/

# EUROPLANET VESPA (EPN-TAP)

# Big Data handling

VO Space     Moving big tables across (load only results)

SSO     Authentication, authorization, groups and consortia

UWS     Universal worker service (job synch, asynch)

SIM-DB     Simulations, theory data

Science platforms   for BD analysis and ML

(*SciServer JHU, NOAO DataLab, CANFAR, Gaia, Euclid , Rubin x Pangeo* )

# VO  in IAU

# Tutorials of VO

www.ivoa.net      portal

https://hendhd.github.io/ivoa_newcomers/

https://www.canfar.net/storage/vault/list/IVOA/virtual2021a  (video)

IVOA Interoperability meetings (May + November)
Newcomers  Intro

Number of VO Schools

EURO-VO DCA, AIDA, ICE
CoSADIE, ASTERICS, ESCAPE

# *Astroinformatics*

# Data-Knowledge-Wisdom Pyramid

# X-informatics



Changing methodology of the Science

Synergy between different worlds

Sociological aspects (net-based research communities)

## Experimental astronomy has become a three players game



- **astronomy**: problems, data, understanding of the data structure and biases

- **mathematics**: evaluation of the data, falsification/validation of theories/models, etc

- **computer science**: implementation of infrastructures, databases, middleware, scalable tools, etc

- Astroinformatics: AAS n. 215, Washington, December 2009, chairperson: K. Borne
- Astroinformatics 2010: Caltech (USA) June 16-19 2010;co-chairpersons: S.G. Djorgovski, G. Longo
- Astroinformatics 2011: UNINA – Sorrento, co-chairpersons: S.G. Djorgovski, G. Longo

Longo 2010

**Need for a new science: Astroinformatics**
*Knowledge Discovery in Databases*

Data Gathering (e.g., from sensor networks, telescopes…)

↳ Data Farming:
    Storage/Archiving
    Indexing, Searchability
    Data Fusion, Interoperability, ontologies, etc.

← Database technologies

Data Mining (or Knowledge Discovery in Databases):
    Pattern or correlation search
    Clustering analysis, automated classification
    Outlier / anomaly searches
    Hyperdimensional visualization

← Key mathematical issues

Data understanding
    Computer aided understanding
    KDD
    Etc.

← Ongoing research

New Knowledge

Longo 2009

Data Mining is the activity of extracting **USEFUL** information from **COMPLEX** data using Statistical Pattern Recognition and Machine Learning methods.

**DM Taxonomy**



1. To catalogue the known (classification)

2. Characterize the unknown (clustering)

3. Find functional dependencies (regression)

4. Find exceptions (outliers)

Supervised Methods

Patterns are learnt from extensive set of templates (Base of Knowledge = BoK)

Unsupervised Methods

Patterns are discovered using the data themselves

# New e-Science Collaborations

## Center for Data-Driven Discovery

- A new research center at Caltech
  - Serves research efforts Institute-wide

- A part of a new, Caltech-JPL joint initiative for data science and technology

- The goals are to assist faculty in **formulation and execution of data-intensive projects**, and facilitate **interdisciplinary sharing of methods, ideas**, novel projects, etc.

Djorgovski

# Data Driven Science

## What is Fundamentally New Here?

- The *information volumes and rates* grow exponentially

  ⟹ *Most data will never be seen by humans*

- A great increase in the data *information content*

  ⟹ *Data driven vs. hypothesis driven science*

- A great increase in the *information complexity*

  ⟹ *There are patterns in the data that cannot be comprehended by humans directly*

Djorgovski

# Hidden Patterns in Data



Pattern or structure (Correlations, Clustering, Outliers, etc.) Discovery in High-Dimensional Parameter Spaces

D >> 3 parameter space hypercube

High-D data cloud: mostly noise, of an arbitrary distribution

But in some corner of some sub-D projection of this data space, there is *something ≠ noise*

Djorgovski

# Visualization in Machine Learning



## A Key Challenge: Visualisating Multidimensional Data Spaces

- Hyperdimensional structures (clusters, correlations, etc.) may be present in many complex data sets, whose dimensionality may be $D \sim 10^2 - 10^4$, or higher

- It is a matter of *data understanding*, choosing the right data mining algorithms, and interpreting the results

- We are biologically limited to perceiving up to $\sim 3 - 12(?)$ dimensions

**What good are the data if we cannot effectively extract knowledge from them?**

Djorgovski

# Visualization of 1 B points – Gaia DR1

# Visualization of Big Data

# Star Forming Regions in Galaxy



Hi-GAL
the Herschel infrared Galactic Plane Survey

70-160-250μm composite

from cold starless clumps to hot HII Regions

Sergio Molinari, INAF-IAPS          IAU Astroinformatics 2016, Sorrento
Credits: Gianluca Li Causi (INAF-IAPS)                    Molinari et al. 2016

# CAVE2 Monash University AU



8m diameter, 330 deg FOV , 80x LCD 46"  1366x768 Stereo + head tracking …...

http://eas.unige.ch/EWASS2017/session.jsp?id=S14

# Astroinformatics in IAU



INTERNATIONAL ASTRONOMICAL UNION

Home | About IAU | IAU Values | Donate | Member Directory | Site Map | Contac

News | Science | Publications | Administration | Training in Astronomy | Astronomy for Education

Home / Science / Scientific Bodies / Commissions / Commission B3 Structure » Commission B3 Astroinformatics and Astrostatistics

## B3 – Commission B3 Astroinformatics and Astrostatistics

### Description

In most of the 20th century, astronomers investigated cosmic phenomena by careful study of individual objects or small samples of planets, stars, galaxies and diffuse media. Datasets were often modest in size with zero (photometry), one (spectra, light curves), or two (images) dimensions. But in the 21st century, increasing resources are devoted to wide-field astronomical surveys, three- or multi-dimensional data, and high-throughput instruments that produce peta-scale datasets and giga-scale samples. In addition to the growing tasks of data reduction, science analysis is becoming more complex. Astronomical insights require characterizing structure in images, spectra or time series. Astrophysical insights require fitting nonlinear, sometimes high-dimensional models to data. Modeling involves both small and large datasets.

IAU Commission B3 focuses on the statistical, computational methodological challenges arising in the various fields of astronomy. It assists the astronomical community in learning existing, and developing new, advanced methodologies to accomplish its goals in this changing environment. The Commission encourages liaison with professional communities in the fields of statistics, applied mathematics and computer science, and with private enterprises. It sponsors meetings and discussions to promulgate advanced methodologies to seek the best scientific insights from the growing flow of data.

Commission Web Page

Commission Members (270)

---

INTERNATIONAL ASTRONOMICAL UNION

Home | About IAU | IAU Values | Donate | Member Directory | Site Map | Contact Us | Login

News | Science | Publications | Administration | Training in Astronomy | Astronomy for Education | Astronomy for Development | Astronomy for the Public

Home / Science / Scientific Bodies / Commissions / Commission B3 Structure / Commission B3 Homepage

## Commission B3 Astroinformatics and Astrostatistics

Scientific Objectives | Members | News | Meetings | Documents | Useful Resources



The unsupervised morphological classification of 200 000 Radio-Galaxy-Zoo images using self-organizing Kohonen map (Polsterer, Gieseke, Igel, 2015ASPC..495...81P). Credit: Kai Lars Polsterer.

Please send any updated email addresses or changes of Institutes to the IAU Secretariat at: iauinfos@iap.fr .

# Remote sensing – Big Data
# Machine Learning



Precise farming

Forestry

Ore mining

Water resources monitoring

Automatic classification of terrain

Resistence of buildings (Aquilla)

# Visualization of Big Data



Breddels 2016

# AstroGeoInformatics

# Astro-Neurology

Description: Detecting objects from astronomical measurements by evaluating light measurements in pixels using intelligent software algorithms.
Image Credit: Catalina Sky Survey (CSS), of the Lunar and Planetary Laboratory, University of Arizona, and Catalina Realtime Transient Survey (CRTS), Center for Data-Driven Discovery, Caltech.

# Finding Cancer Signatures NASA



Description: Detecting objects from oncology images using intelligent software algorithms transferred to and from space science.
Image Credit: EDRN Lung Specimen Pathology image example, University of Colorado

# Science Platforms - SciServer

# Digital Pathology based on SDSS

# Digital Pathology based on SDSS



Astronomy viewer

Pathology viewer

# *Artificial Intelligence*

# Human – Machine Collaboration

From the Information Technology to the
Cognition Technology: Towards a
Human-Computer Collaborative Discovery

AS WE MAY THINK
A TOP U.S. SCIENTIST FORESEES A POSSIBLE FUTURE WORLI
IN WHICH MAN-MADE MACHINES WILL START TO THINI
by VANNEVAR BUSH

Vannevar Bush
(1945)

Man-Computer    Symbiosis

J.C.R. Licklider
(1960)

MEMEX

„Hypertext"
Associations

Global network
Interactivity
Graphics

# Turing Test

# MIND

## A QUARTERLY REVIEW

OF

## PSYCHOLOGY AND PHILOSOPHY

❦

## I.—COMPUTING MACHINERY AND INTELLIGENCE

### BY A. M. TURING

1.   *The Imitation Game.*

# Imitation Game

The new form of the problem can be described in terms of a game which we call the 'imitation game'. It is played with three people, a man (A), a woman (B), and an interrogator (C) who may be of either sex. The interrogator stays in a room apart from the other two. The object of the game for the interrogator is to determine which of the other two is the man and which is the woman. He knows them by labels X and Y, and at the end

Now suppose X is actually A, then A must answer. It is A's object in the game to try and cause C to make the wrong identification. His answer might therefore be

'My hair is shingled, and the longest strands are about nine inches long.'

question and answers can be repeated by an intermediary. The object of the game for the third player (B) is to help the interrogator. The best strategy for her is probably to give truthful answers. She can add such things as 'I am the woman, don't listen to him!' to her answers, but it will avail nothing as the man can make similar remarks.

# Turing Test – Core Idea

We now ask the question, 'What will happen when a machine takes the part of A in this game?' Will the interrogator decide wrongly as often when the game is played like this as he does when the game is played between a man and a woman? These questions replace our original, 'Can machines think?'

Strong  AI  (passes TT)   x     Weak  AI
Searle 1980 – Chinese Room

Turing prize   ~   Nobel prize in IT
1954 suicide by apple with cyanide

# Origins of term „Artificial Intelligence"





1956 Dartmouth Conference: The Founding Fathers of AI

Founding fathers of AI. Courtesy of scienceabc.com

# First Chatbot - ELIZA

Weizenbaum  MIT   1964-67



```
Did you come to me because you are not ?

>i dont know
What does that suggest to you?

>that i am unintelligent
How long have you been unintelligent?

>since i was born
I see.

>arent you supposed to cheer me up
You're not really talking about me, are you?

>no im talking about me
Why not?

>mrh
I'm not sure I understand you fully.

>rg
Come come elucidate your thoughts.

>what is elucidate
What do you think?

>to see
Can you elaborate on that?

>
```

```
Welcome to
            EEEEEE  LL       IIII   ZZZZZZ   AAAAA
            EE      LL        II        ZZ  AA   AA
            EEEEE   LL        II       ZZZ  AAAAAAA
            EE      LL        II      ZZ    AA   AA
            EEEEEE  LLLLLL  IIII  ZZZZZZ    AA   AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landsteiner 2005.


ELIZA: Is something troubling you ?
YOU:   Men are all alike.
ELIZA: What is the connection, do you suppose ?
YOU:   They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU:   Well, my boyfriend made me come here.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU:   He says I'm depressed much of the time.
ELIZA: I am sorry to hear that you are depressed.
YOU:   It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:
```

# Recent Advanced Neural Networks

Reinforcement Learning

Active learning (oracle)

Bayesian Deep Learning

Physics Informed (Aware) NN
Theory is important again

GANs (image creation)

Diffusion models

Foundation models



**NEURAL NETWORK ARCHITECTURE TYPES**

SINGLE LAYER PERCEPTRON

RADIAL BASIS NETWORK

MULTI LAYER PERCEPTRON

RECURRENT NEURAL NETWORK

LSTM RECURRENT NEURAL NETWORK

HOPFIELD NETWORK

BOLTZMANN MACHINE

INPUT UNIT    HIDDEN UNIT    BACKFED INPUT UNIT

OUTPUT UNIT    FEEDBACK WITH MEMORY UNIT    PROBABILISTIC HIDDEN UNIT

# Reviews of ML/AI in Astronomy

DRAFT VERSION APRIL 17, 2019
Typeset using LaTeX preprint2 style in AASTeX61

MACHINE LEARNING IN ASTRONOMY: A PRACTICAL OVERVIEW

DALYA BARON[1]

[1]School of Physics and Astronomy
Tel-Aviv University
Tel Aviv 69978, Israel

To appear in: *Artificial Intelligence for Science*,
eds. A. Choudhary, G. Fox and T. Hey
Singapore: World Scientific, in press (2023)

## Applications of AI in Astronomy

S. G. Djorgovski*, A. A. Mahabal*, M. J. Graham*, K. Polsterer[†],
A. Krone-Martins[‡]

Experimental Astronomy (2022) 53:1–43
https://doi.org/10.1007/s10686-021-09827-4

**REVIEW ARTICLE**

Check for updates

## Astronomical big data processing using machine learning: A comprehensive review

Snigdha Sen[1,2] · Sonali Agarwal[1] · Pavan Chakraborty[1] ·
Krishna Pratap Singh[1]

Received: 15 July 2021 / Accepted: 27 December 2021 / Published online: 14 January 2022
© The Author(s), under exclusive licence to Springer Nature B.V. 2022

Astronomy and Computing 48 (2024) 100851

Contents lists available at ScienceDirect

## Astronomy and Computing

ELSEVIER

journal homepage: www.elsevier.com/locate/ascom

## A review of unsupervised learning in astronomy

S. Fotopoulou

School of Physics, HH Wills Physics Laboratory, University of Bristol, Tyndall Avenue, Bristol, BS8 1TL, United Kingdom

# Idea of LLMs

Tomáš Mikolov

Dissertation FIT VUT Brno 2012,
MS Research, Google Brain, Facebook, CIIRC

word2vec     sentence as a vector in big-N space
            improves Google Translate

Fasttext.cc

https://www.seznamzpravy.cz/clanek/jeho-figl-v-usa-nakopl-vyvoj-strojoveho-uceni-proc-se-expert-vratil-176156

# Origin of the Current AI Boom

## Attention Is All You Need

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

**Aidan N. Gomez**[*][†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[*][‡]
illia.polosukhin@gmail.com

# Attention



She is eating a green apple

Low attention

High attention

$$\text{attention}(Q, K, V) = \text{softmax}\left(\frac{QK^{T}}{\sqrt{d_k}}\right) V$$

"The cat drank the milk because it was hungry"

"The cat drank the milk because it was sweet"

| The | | The | |
| cat | | cat | |
| drank | | drank | |
| the | | the | |
| milk | | milk | |
| because | | because | |
| it | | it | |
| was | | was | |
| hungry | | hungry | |

| The | | The | |
| cat | | cat | |
| drank | | drank | |
| the | | the | |
| milk | | milk | |
| because | | because | |
| it | | it | |
| was | | was | |
| sweet | | sweet | |

K. Doshi

# Transformer
# (GPT=Generative Pre-trained Transformer)

# Transformers in Astronomy

- ASTROMER: A transformer-based embedding for the representation of light curves
  - pre-trained on millions of light curves from different surveys (MACHO, OGLE, ATLAS)
  - representation to create informative light curves embeddings
  - finetuned for solving downstream tasks, e.g. classification of variable stars, predicting physical parameters

**AST ROMER**

https://www.stellardnn.org/projects/astromer/index.html

C. Donoso-Oliva et al. ASTROMER: A transformer-based embedding for the representation of light curves.

# LLMs in 2024

## BILLBOARD CHART FOR LANGUAGE MODELS — JUN/2024

### MODELS

| Now (Jun/2024) | 6m ago (Dec/2023) | 12m ago (Jun/2023) | ALScore | Model name / Details | AI lab / Openness |
|---|---|---|---|---|---|
| ❶ | — | — | 29.8 | **Claude 3 Opus** 2T trained on 40T tokens* | ◆ Anthropic API |
| ❷ | 1 | — | 22.4 | **Gemini Ultra 1.0** 1.5T trained on 30T tokens* | ◆ Google DM API |
| ❸ | — | — | 22.4 | **Gemini 1.5 Pro** 1.5T trained on 30T tokens* | ◆ Google DM API |
| ❹ | — | — | 21.1 | **Yi-XLarge** 2T trained on 20T tokens* | ◆ 01-ai API |
| ❺ | — | — | 16.3 | **Inflection-2.5** 1.2T on 20T tokens* | ◆ Inflection AI API |
| ❻ | 2 | 1 | 15.9 | **GPT-4 (family)** 1.7T trained on 13T tokens* | ◆ OpenAI API |
| ❼ | 3 | — | 14.9 | **ERNIE 4.0** 1T trained on 20T tokens* | ◆ Baidu API |
| ❽ | — | — | 8.2 | **SenseNova 5.0** 600B on 10T tokens | ◆ SenseTime API |

### DATASETS

| Now (Jun/2024) | 6m ago (Dec/2023) | 12m ago (Jun/2023) | Size (TB) | Dataset name / Details | AI lab / Language |
|---|---|---|---|---|---|
| ❶ | 1 | — | 130 | **Gemini** 30T tokens in 130TB* | ◆ Google DM Multilingual |
| ❷ | 2 | — | 125 | **RedPajama-Data-v2** 30T tokens in 125TB | ◆ Together AI Multilingual |
| ❸ | 3 | 1 | 86 | **Piper monorepo** 37.9T tokens in 86TB | ◆ Google Code |
| ❹ | 4 | — | 40 | **Massive Never-ending BT Vast Chinese corpus** 30T/40TB | ◆ MNBVC Chinese |
| ❺ | — | — | 44 | **FineWeb** 15T tokens in 44TB | ◆ HF English |
| ❻ | 5 | 2 | 40 | **GPT-4** 13T tokens in 40TB* | ◆ OpenAI English |
| ❼ | — | — | 31.5 | **FineWeb-Edu-score-2** 5.4T tokens in 31.5TB | ◆ HF English |
| ❽ | 6 | — | 27 | **CulturaX** 6.3T tokens in 27TB | ◆ UOregon Multilingual |

Selected highlights only, some older models disregarded. * = estimates and hypothesis only based on current information. Alan D. Thompson. June 2024. https://lifearchitect.ai/

🔗 LifeArchitect.ai/models

# LLMs in Astronomy



ESO Cosmic Duologues
https://www.youtube.com/@ESOCosmicDuologues/featured

AstroLLaMA : 7B pararameters,
300 000 abstracts from ADS

Cosmology – FIT SED

ESO User Man

Proposal rewiews

Project assessments

Ethical problems

# Review of LLMs in Astronomy

## Astronomia ex machina: a history, primer and outlook on neural networks in astronomy

Michael J. Smith and James E. Geach

Department of Physics, Astronomy and Mathematics, School of Physics, Engineering and Computer Science, University of Hertfordshire, Hatfield AL10 9AB, UK

MJS, 0000-0003-0220-5125; JEG, 0000-0003-4964-4635

In this review, we explore the historical development and future prospects of artificial intelligence (AI) and deep learning in astronomy. We trace the evolution of connectionism in astronomy through its three waves, from the early use of multilayer perceptrons, to the rise of convolutional and recurrent neural networks, and finally to the current era of unsupervised and generative deep learning methods. With the exponential growth of astronomical data, deep learning techniques offer an unprecedented opportunity to uncover valuable insights and tackle previously intractable problems. As we enter the anticipated fourth wave of astronomical connectionism, we argue for the adoption of GPT-like foundation models fine-tuned for astronomical applications. Such models could harness the wealth of high-quality, multimodal astronomical data to serve state-of-the-art downstream tasks. To keep pace with advancements driven by Big Tech, we propose a collaborative, open-source approach within the astronomy community to develop and maintain these foundation models, fostering a symbiotic relationship between AI and astronomy that capitalizes on the unique strengths of both fields.

## Designing an Evaluation Framework for Large Language Models in Astronomy Research

| John F. Wu | Alina Hyk | Kiera McCormick |
| Christine Ye | Simone Astarita | Elina Baral |
| Jo Ciuca | Jesse Cranney | Anjalie Field |
| Kartheik Iyer | Philipp Koehn | Jenn Kotler |
| Sandor Kruk | Michelle Ntampaka | Charles O'Neill |
| Joshua E.G. Peek | Sanjib Sharma | Mikaeel Yunus |

## Large Language Models: A Survey

Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu
Richard Socher, Xavier Amatriain, Jianfeng Gao
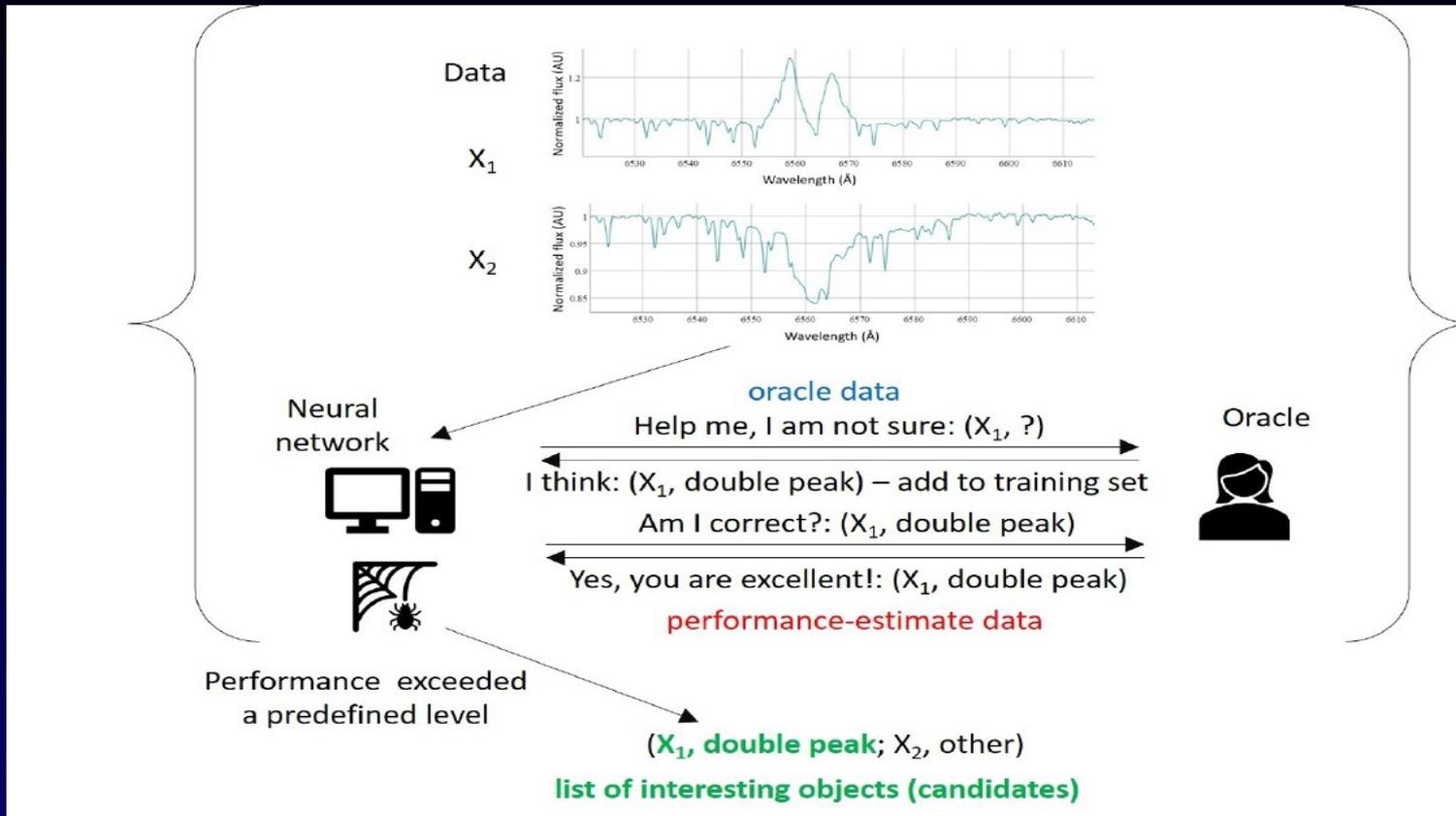
# Retrival Augmented Generation (RAG)



https://arxiv.org/pdf/2405.20389

# Conclusions

AI already helps us to discover the secrets of the  Universe

But the future is ..... uncertain !

TRUST   RELIABILITY  FAKES

# Backup Slides

# Active Learning (insufficient labels)



Oracle :  Human – Machine Interaction

# Publication in A&A

## Active deep learning method for the discovery of objects of interest in large spectroscopic surveys[*,**]

P. Škoda[1,2], O. Podsztavek[2] and P. Tvrdík[2]

+

# New Catalogue on Vizier



Active deep learning in large spectros. surveys : J/A+A/643/A122

Access to: VizieR | FTP | ReadMe | TAP | Xmatch

Authors : Skoda P. , Podsztavek O., Tvrdik P.

VizieR DOI : 10.26093/cds/vizier.36430122  Cite
Bibcode : 2020A&A...643A.122S (ADS)

UAT : Emission line stars, Surveys, Spectroscopy

Compilation (CCC)

Inserted into VizieR : 11-Nov-2020
Last modification : 02-Feb-2021

Article Origin | Description | See also | Prov | FTP | VizieR

Active deep learning method for discovery of objects of interest in large spectroscopic surveys. (2020)
Go to the original article (10.1051/0004-6361/201936090)

Keywords : surveys - virtual observatory tools - methods statistical - techniques: spectroscopic - stars: emission-line, Be - line: profiles

Abstract:Current archives of the LAMOST telescope contain millions of pipeline-processed spectra that have probably never been seen by human eyes. Most of the rare objects with interesting physical properties, however, can only be identified by visual analysis of their characteristic spectral features. A proper combination of interactive visualisation with modern machine learning techniques opens new ways to discover such objects. We apply active learning classification methods supported by deep convolutional neural networks to automatically identify complex emission-line shapes in multi-million spectra archives. We used the pool-based uncertainty sampling active learning method driven by a custom-designed deep convolutional neural network with 12 layers. The architecture of the network was inspired by VGGNet, AlexNet, and ZFNet, but it was adapted for operating on one-dimensional feature vectors. The unlabelled pool set is represented by 4.1 million spectra from the LAMOST data release 2 survey. The initial training of the network was performed on a labelled set of about 13000 spectra obtained in the 400Å wide region around Hα by the 2m Perek telescope of the Ondrejov observatory, which mostly contains spectra of Be and related early-type stars. The differences between the Ondrejov intermediate-resolution and the LAMOST low-resolution spectrographs were compensated for by Gaussian blurring and wavelength conversion. After several iterations, the network was able to successfully identify emission-line stars with an error smaller than 6.5%. Using the technology of the Virtual Observatory to visualise the results, we discovered 1013 spectra of 948 new candidates of emission-line objects in addition to 664 spectra of 549 objects that are listed in SIMBAD and 2644 spectra of 2291 objects identified in an earlier paper of a Chinese group led by Wen Hou. The most interesting objects with unusual spectral properties are discussed in detail. (hide)

# Outreach

# Outreach

- Seznam zprávy
- 24 comments
- No hating !
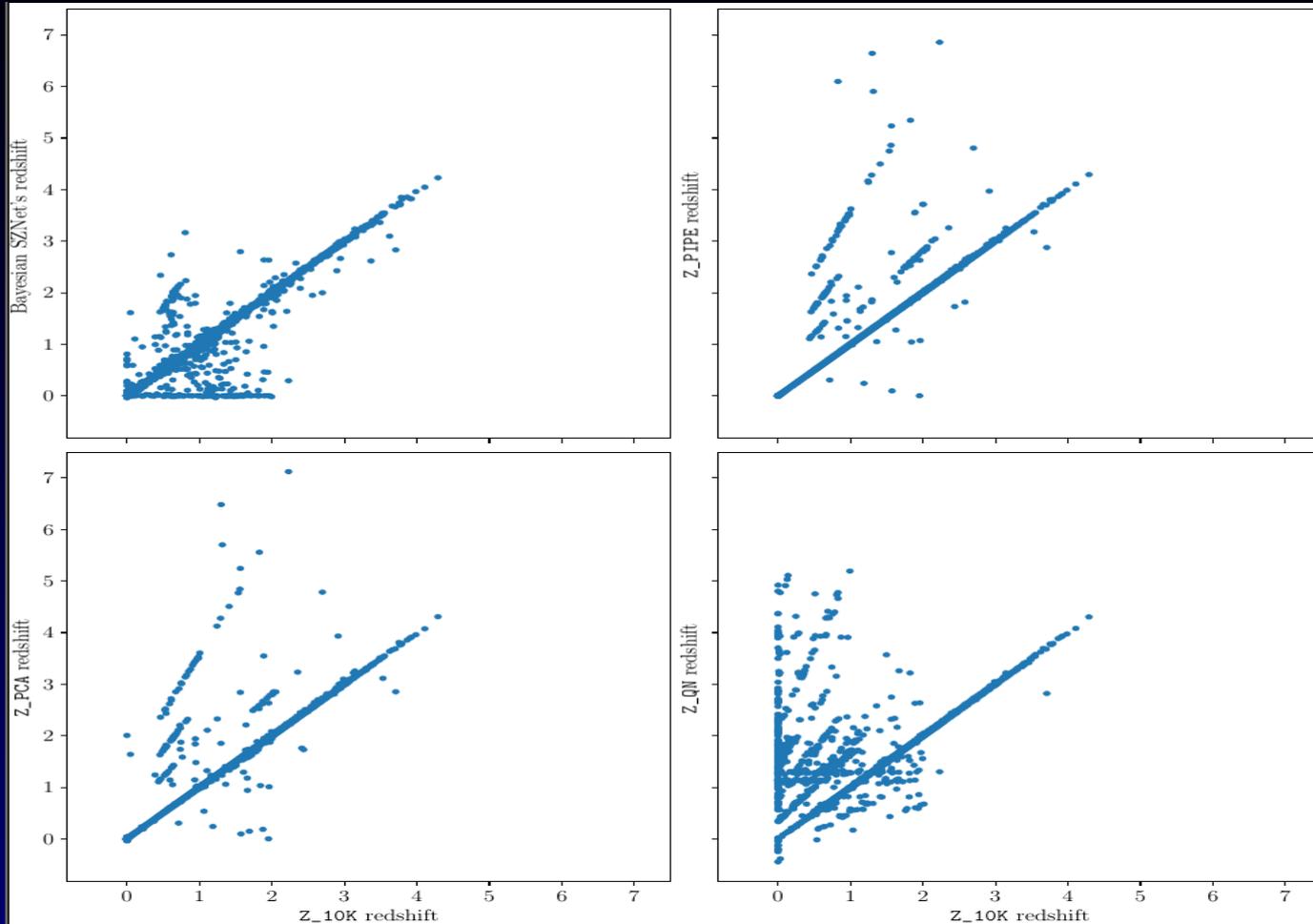
# Redshift



326886 SDSS/BOSS Data Release 12 Quasars

# SDSS Template Library

# Systematic Errors in Pipelines



Z_10K:
~10000 randomly
selected spectra
visually  checked

Used to evaluate SDSS
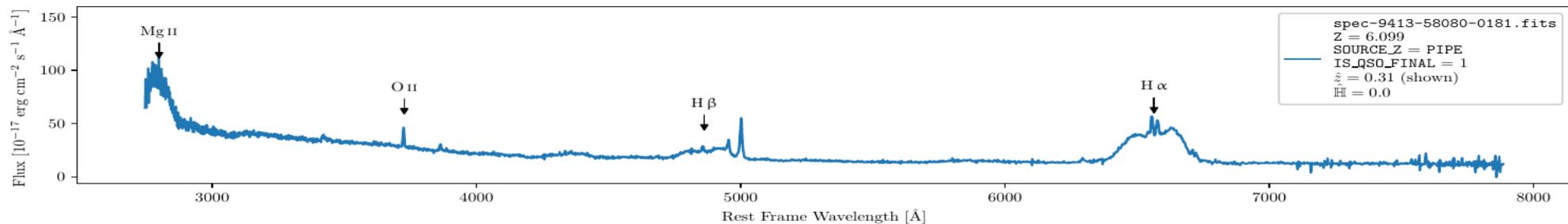DR16Q  pipeline

# BNN corrects the SDSS pipeline



**Figure B6.** Spectrum with incorrectly high redshift prediction by the pipeline. The Bayesian CNN correctly predicted $\hat{z} = 0.31$ with $\hat{\mathbb{H}} = 0$.
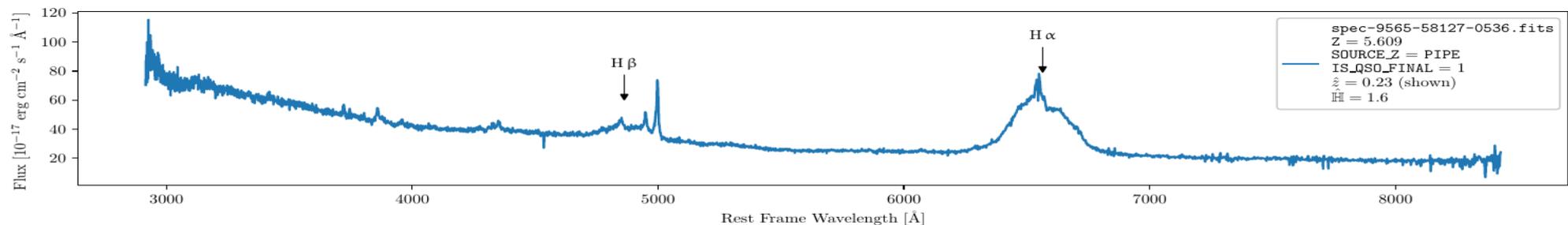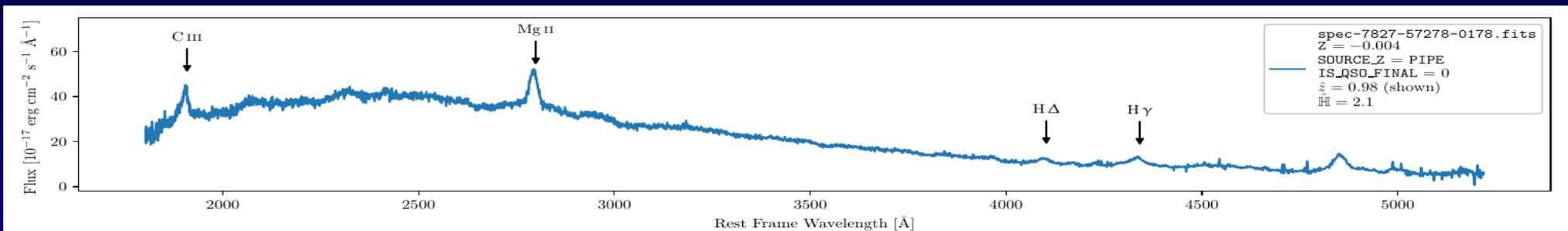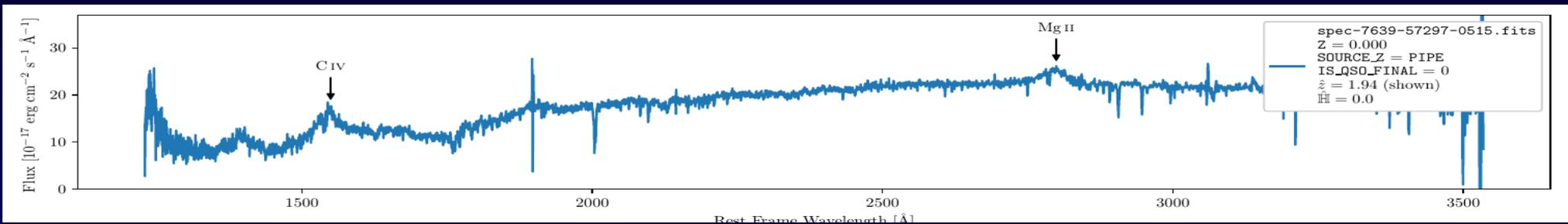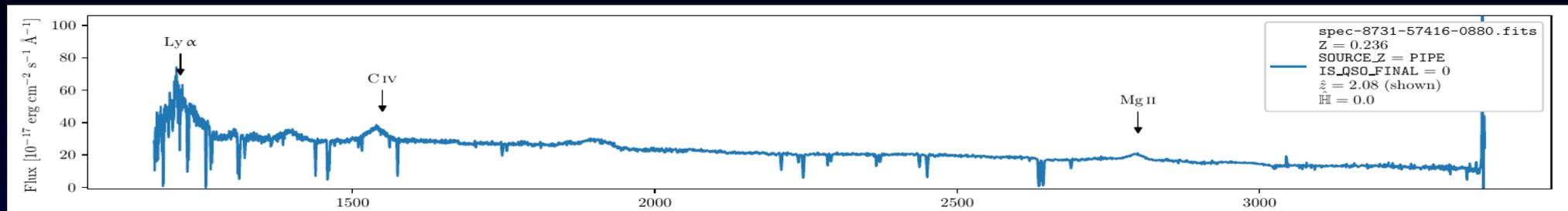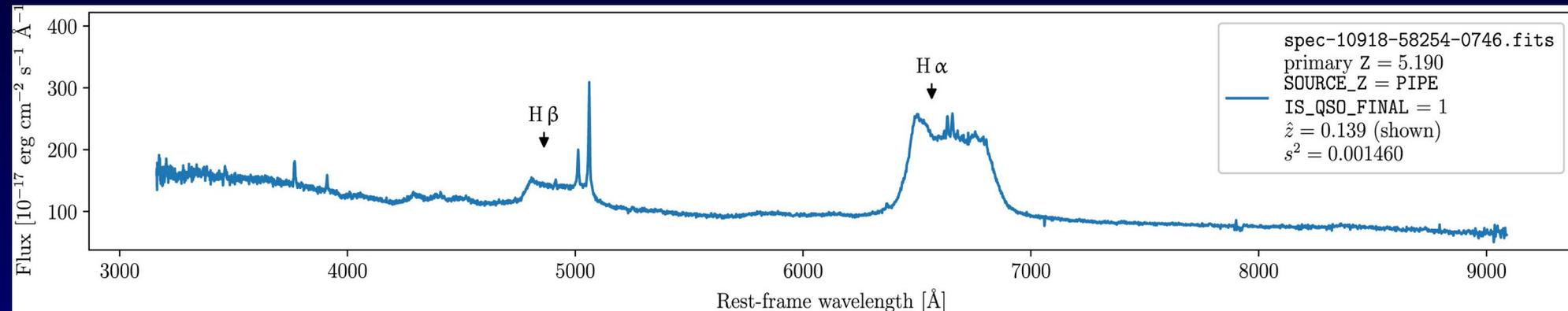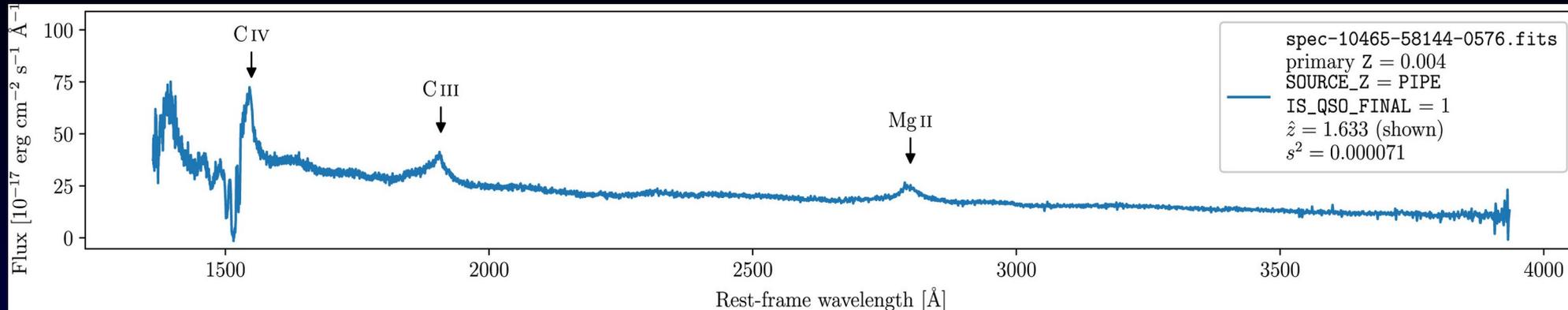


**Figure B7.** Spectrum with incorrectly high redshift prediction by the pipeline. The Bayesian CNN correctly predicted $\hat{z} = 0.23$ with $\hat{\mathbb{H}} = 1.6$.
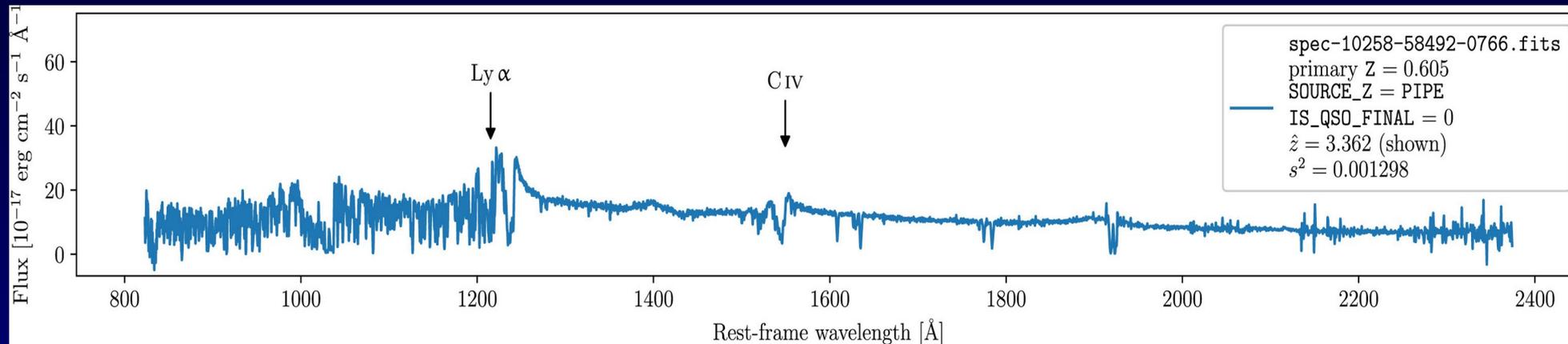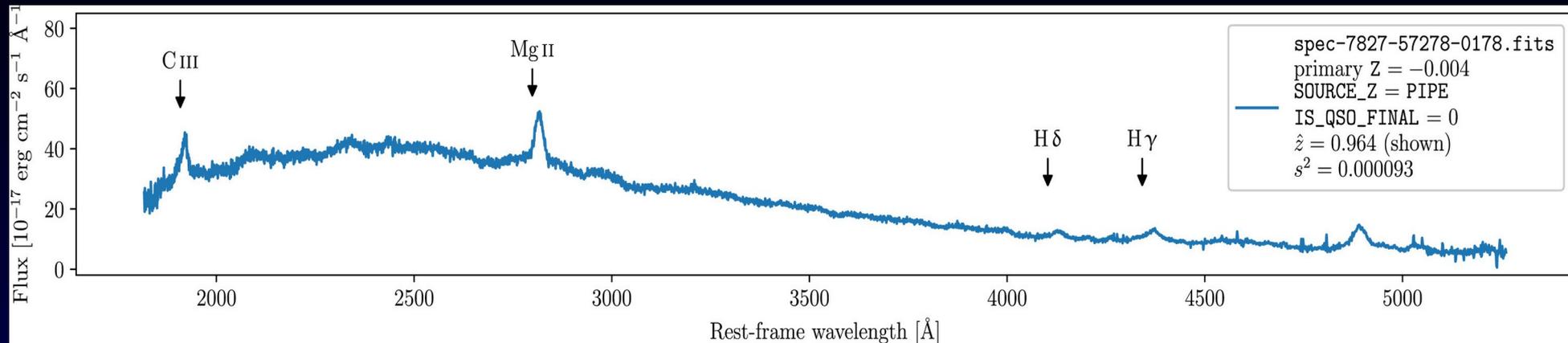
# QSOs missing due to SDSS pipeline error

# BDN Corrects the SDSS Pipeline

# QSOs missing due to SDSS pipeline error

# SDSS Predicts QSO
# But It Is a Star