



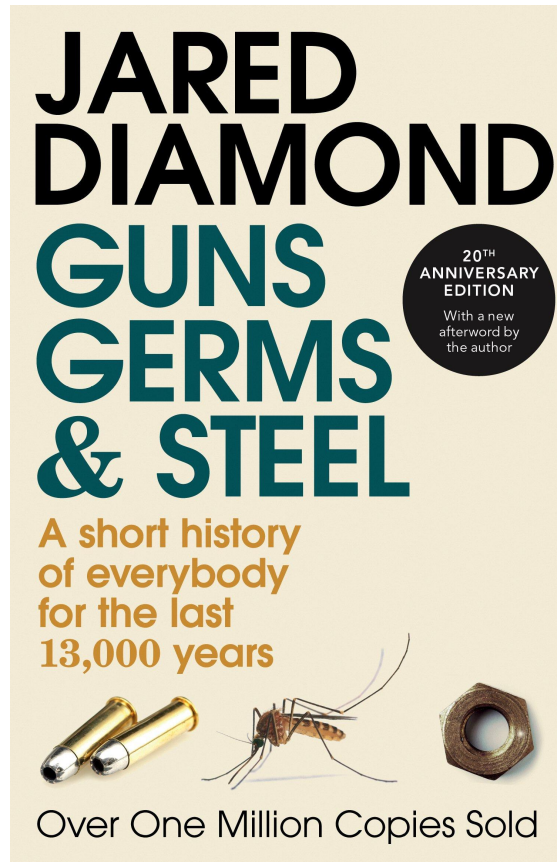
ΧΑΡΟΚΟΠΕΙΟ ΠΑΝΕΠΙΣΤΗΜΙΟ  
HAROKOPIO UNIVERSITY



Vasilis Efthymiou

# BITS STARS & KNOWLEDGE AI at the Frontiers of Astrophysics

Summer500  
Service Oriented Computing





*“What would it take for a Foundation Model that understands (the lifecycle of) our Universe?”*

**Knowledge**

**Bits**

**Stars**





# Not there yet...

June 18, 2025



**CHATGPT-4o** LOSES TO **ATARI**  
FROM **1979** IN BEGINNER-LEVEL  
**CHESS MATCH**



# AI $\neq$ ML

AI





# Outline

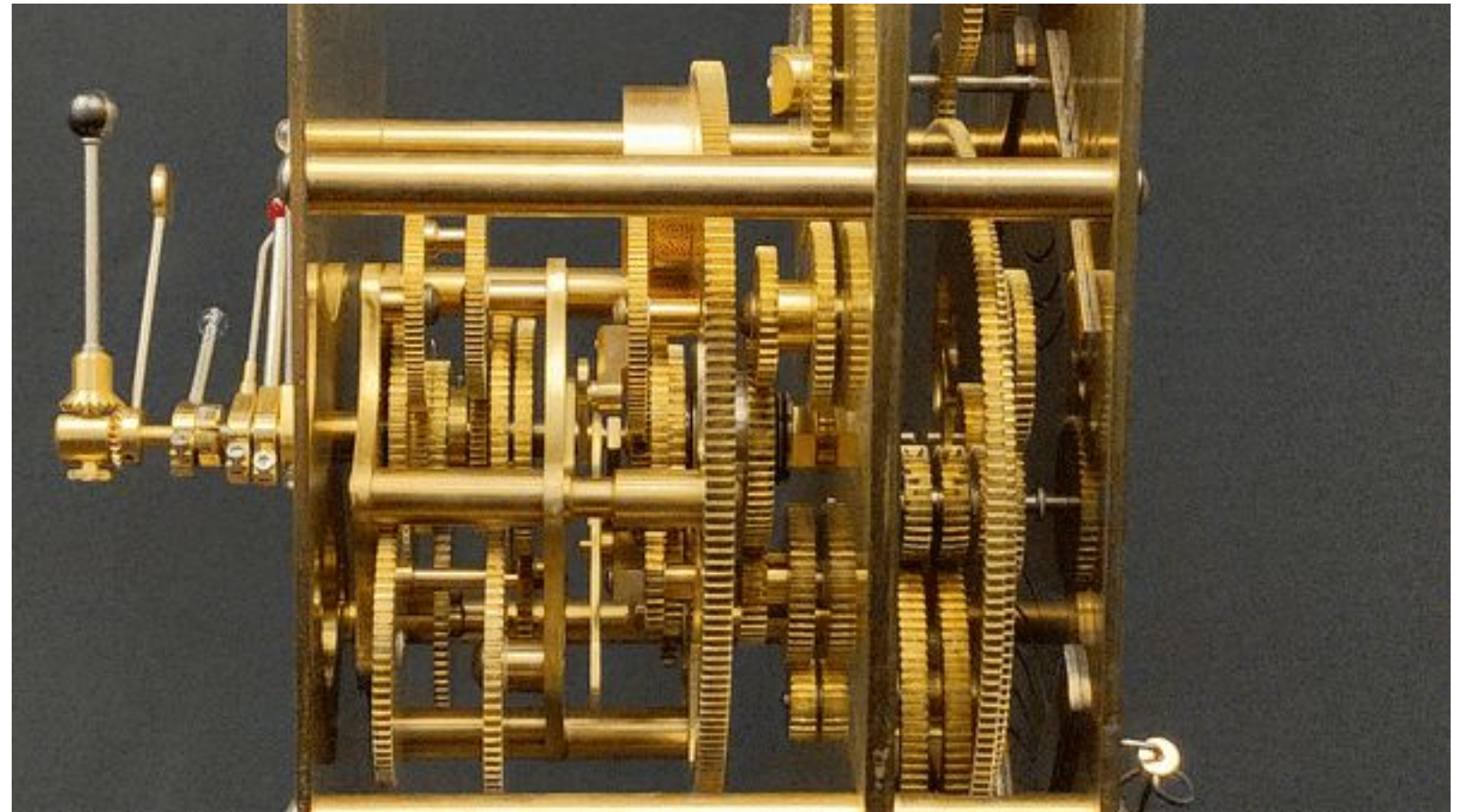
- “Bits” (~ AI)
    - latest advances in deep learning
  - “Stars” (~ astrophysics)
    - challenges in “traditional” astrophysics
- UniversAI: Exploring the Universe with AI
- Knowledge...
    - ... representation (knowledge graphs)
  - The PARSEC project



# Bits



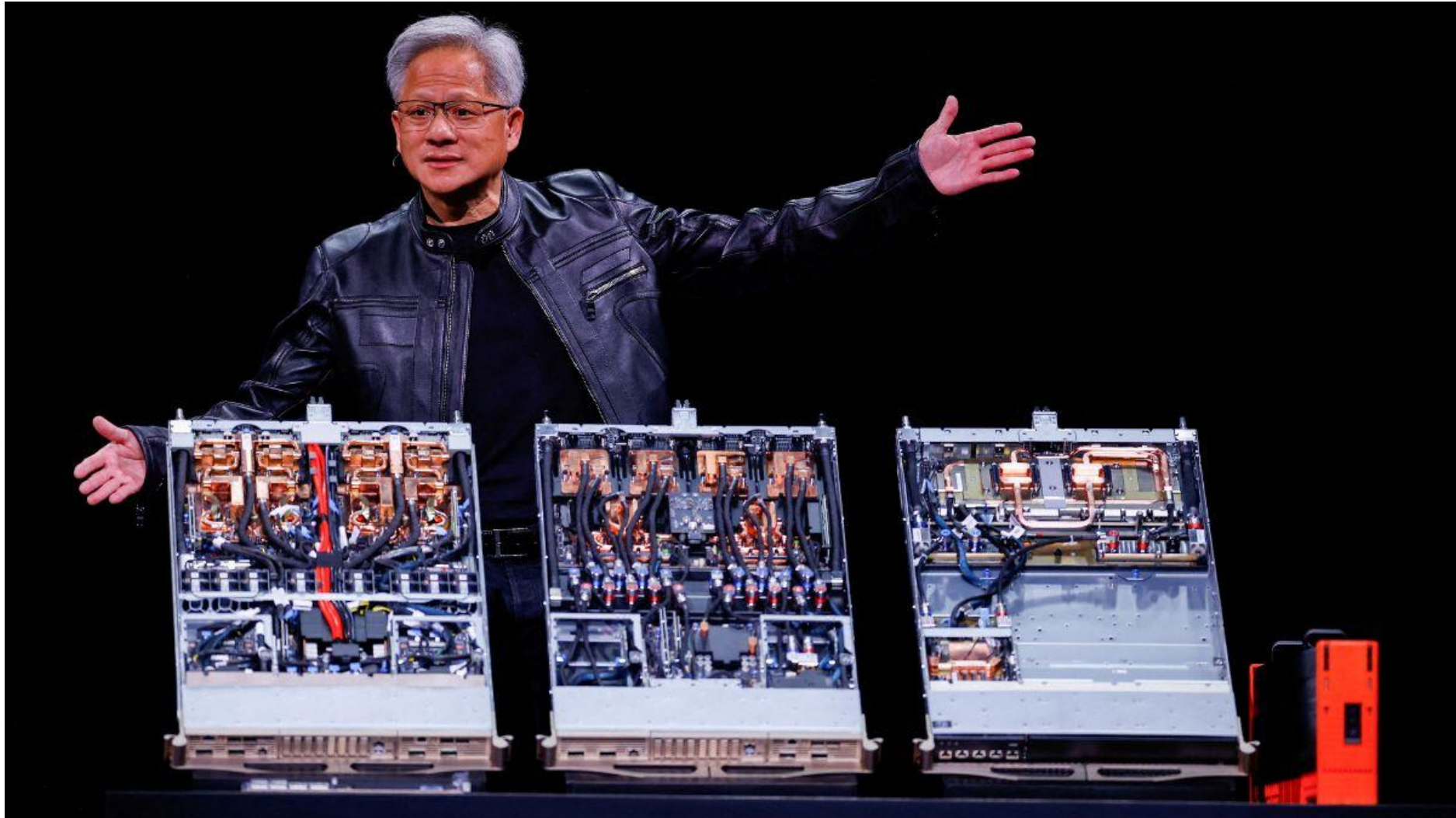
June 18, 2025



SummerSoC 2025, Hersonissos, Crete, Greece



# ... more bits





# Recent advances in AI

- Foundation models for science
  - SciBERT / Galactica / Claude
  - Retrieval-augmented generation (RAG) for data (e.g., catalog) querying
- Generative models for simulations
  - Diffusion models / GANs
- Graph Neural Networks / Graph Transformers
  - Detecting matches/missing links across one or more data sources
- Counterfactuals
  - Understanding how AI models (would have) work(ed)

# Stars





... more stars

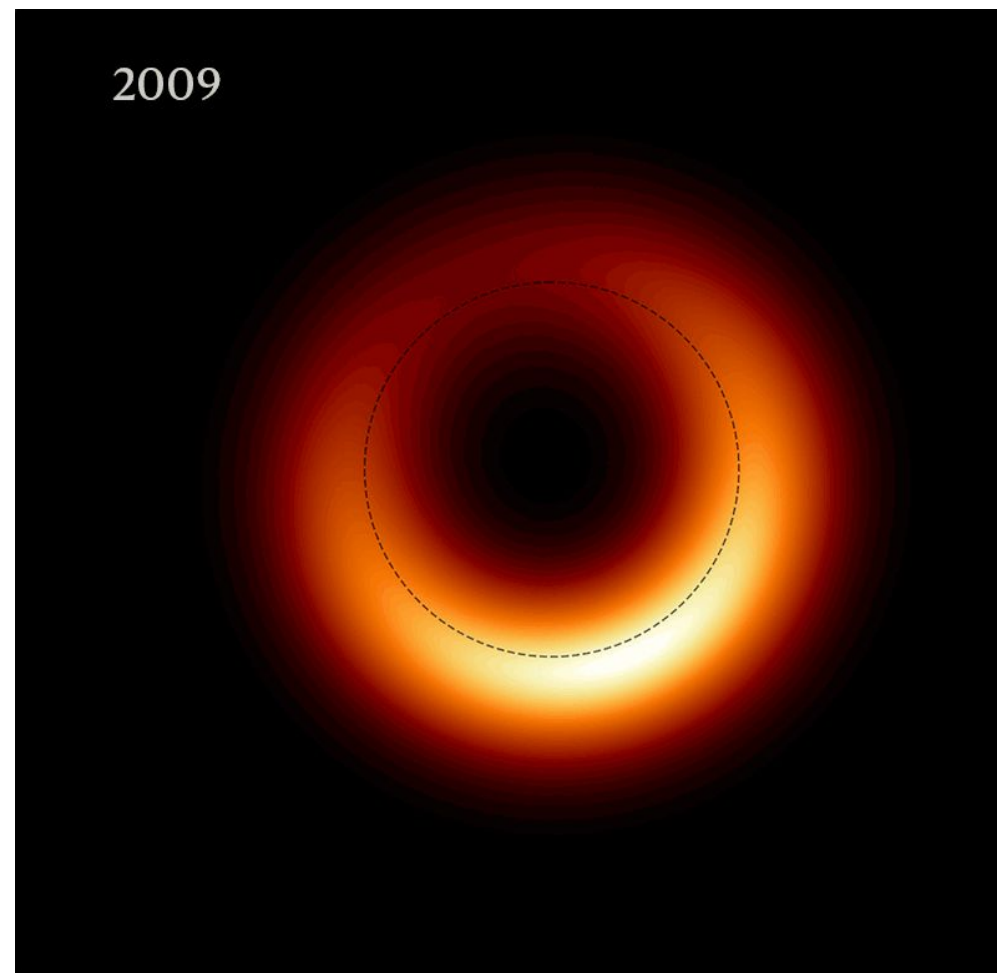
will start in 2025  
20 TB per night  
60 PB in total



# AI & Astrophysics



The first-ever image of a black hole is now a movie...



<https://www.nature.com/articles/d41586-020-02717-3>



# Recent advances in astronomy

- JWST & Multi-wavelength data fusion
  - Heterogeneous, high-res datasets across wavelengths and missions
  - Opportunity: how do we align, interpret and reason about them?
- Astroinformatics
  - Cosmological simulations now produce petabytes of synthetic universes
  - Opportunity: indexing, querying, similarity computations
- Citizen science
  - Galaxy zoo: large crowdsourced labels (*also: largely inconsistent*)
  - Opportunity: weakly/semi-supervised learning

# UniversAI : Exploring the Universe with Artificial Intelligence



June 2-6, 2025, Athens, Greece



# What led to UniversAI?

- Cultural cross-talk
  - What is a “source”?
  - ... “survey” vs “review”
  - ... “referee” vs “reviewer”
  - ... conference vs journal papers
- Falling behind the Big (AI) Bang / space missions / data repositories
- Idea: Let’s bring the best of the two communities together

# Invited talks

- **Jean-Luc Starck** (*CEA –Saclay*)
- **Christos Diou** (*HUA*)
- **Torsten Enßlin** (*MPA*)
- **Themis Palpanas** (*French University Institute*)
- **Tyson Littenberg** (*NASA*)
- **Angela Bonifati** (*Lyon 1 University*)
- **Meghyn Bienvenu** (*CNRS*)
- **Laurent Eyer** (*Geneva Astronomical Observatory*)
- **Federica Bianco** (*University of Delaware*)
- **Georgia Koutrika** (*Athena RC*)
- **Hendrik Müller** (*NRAO*)
- **Grigoris Tsagatakis** (*FORTH*)



# UniversAI – Highlights

## T. Littenberg – GW and AI

The co-evolution of  
gravitational wave astrophysics  
and artificial intelligence



Tyson B. Littenberg @ NASA Marshall Space Flight Center

Noise Reduction

Computational Efficiency

Event Classification

Source Extraction

**For GW detection, interpretability  
is *critical*. Critical.**

# UniversAI – Highlights

## C. Diou - XAI

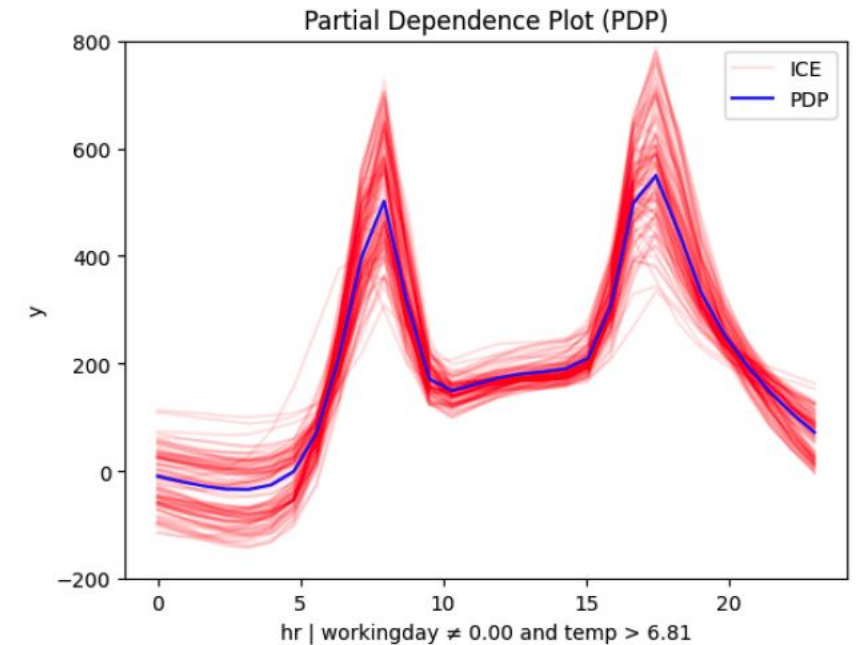
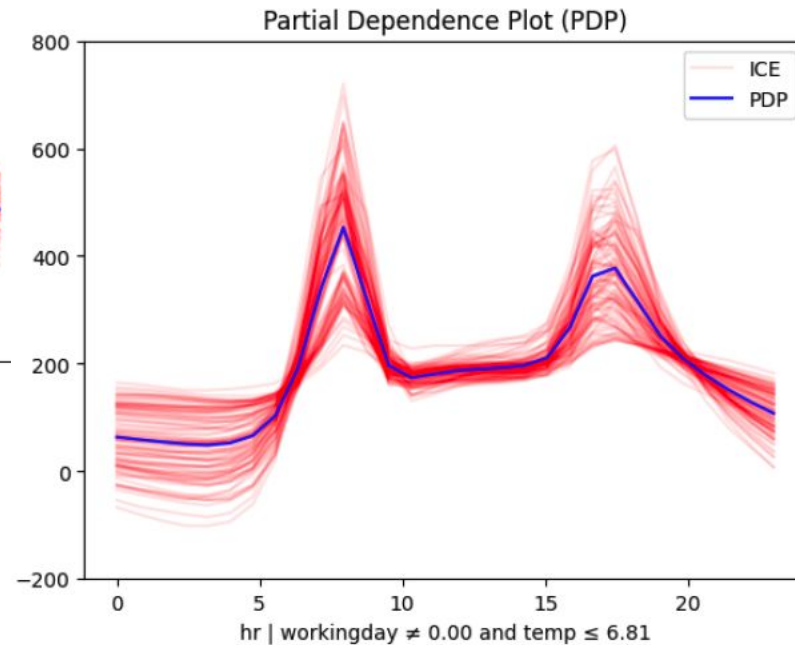
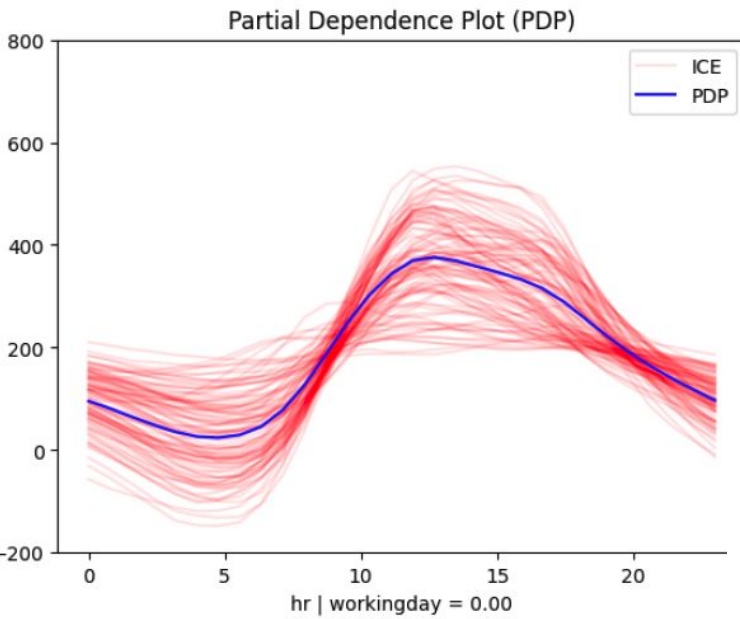
Using regional effect plots for machine learning model  
interpretation: The Effector python package

Christos Diou

Department of Informatics and Telematics  
Harokopio University of Athens

IAUS 397: UniversAI 2025 31/05/2025

1/46





# UniversAI – Highlights

## T. Palpanas – Anomaly Detection

### Machine Learning on Very Large Time Series Collections

#### Similarity Search and Subsequence Anomaly Detection



Themis Palpanas

Université Paris Cité  
French University Institute

diNo 381

## Subsequence Anomaly Detection

Proposed Approach: Series2Graph

**Graph  $G_\ell$  :**  
Given a time series  $T$ , and an input length  $\ell$ , we build a graph  $G_\ell(\mathcal{N}, \mathcal{E})$ , for which:

For instance:

The subsequence  $T_{i,\ell+2}$ , is a path in the graph  
 $P_{th}(T_{i,\ell+2}) = \langle T_{i,\ell}, T_{i+1,\ell}, T_{i+2,\ell} \rangle = \langle N_4, N_5, N_6 \rangle$  in  $G_\ell$ .

The subsequence  $T_{j,\ell+2}$ , is a path in the graph  
 $P_{th}(T_{j,\ell+2}) = \langle T_{j,\ell}, T_{j+1,\ell}, T_{j+2,\ell} \rangle = \langle N_2, N_1, N_2 \rangle$  in  $G_\ell$ .

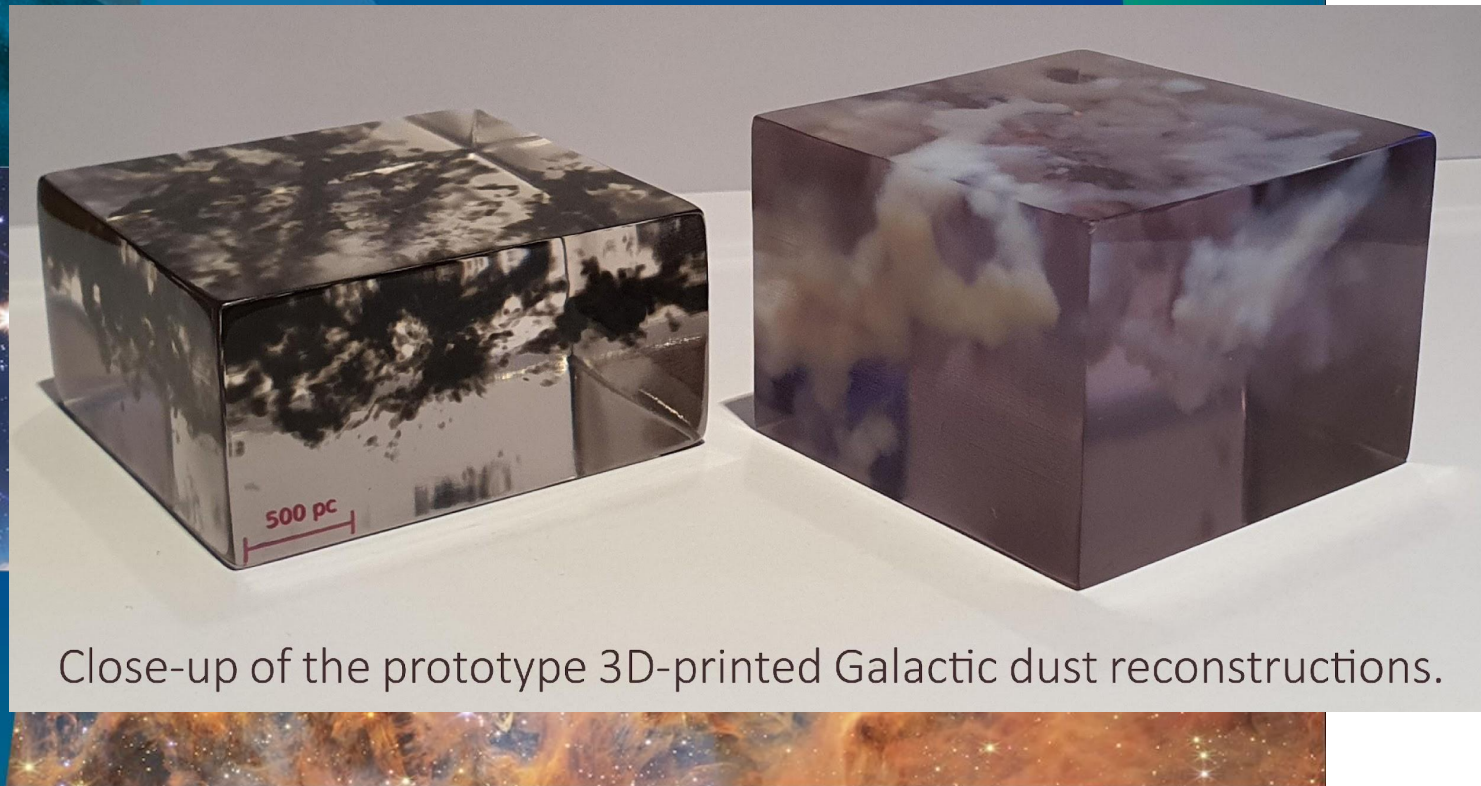
$\text{Norm}(P_{th}(T_{j,\ell+2})) \ll \text{Norm}(P_{th}(T_{i,\ell+2}))$

**unsupervised**

Themis Palpanas - UniversAI - Jun 2025

# UniversAI – Highlights

## T. Enßlin – Milky Way Atlas





# UniversAI – Highlights ... and many more



Vision-Language  
Models for Radio  
Astronomy

Simone Riggi

✉ [simone.riggi@inaf.it](mailto:simone.riggi@inaf.it)



Diffusion Models for Emulating the Large-Scale Structure of the Universe in Modified Gravity Cosmologies

Do Androids Dream of Exploding Stars  
and Receding Galaxies?

E. García-Farieta  
Universidad de Córdoba  
with: JK Riveros, P. Saavedra, H.J. Hortúa, I. Olier

**FALCO:**  
a Foundation model of Astronomical Light Curves  
for time domain astronomy

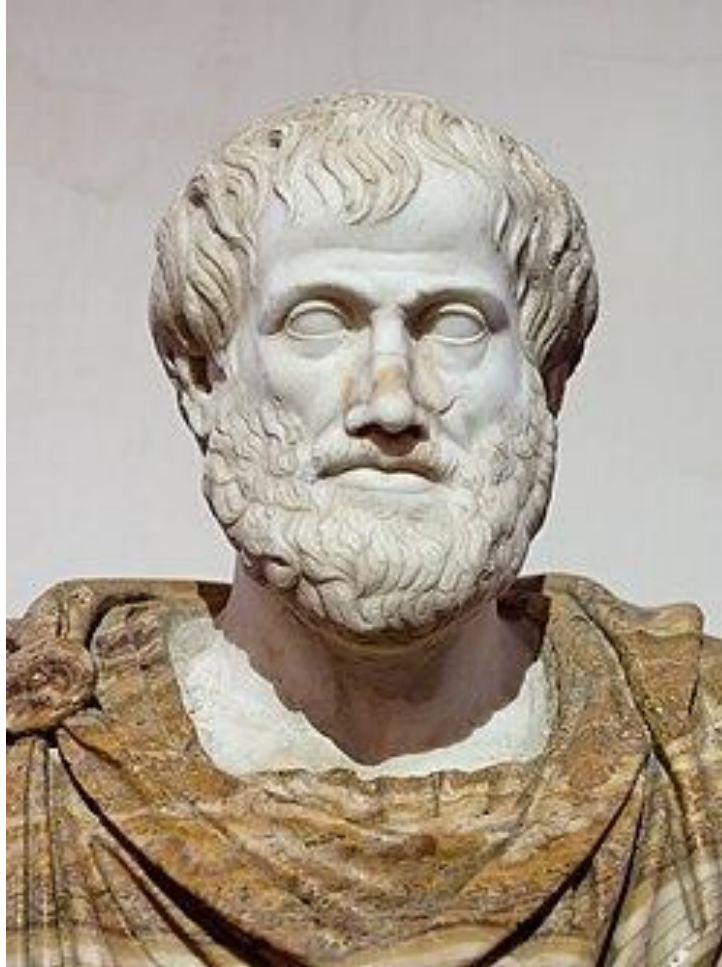
Yihan Tao  
National Astronomical Observatories, Chinese Academy of Sciences  
Xiaoxiong Zuo  
LMU Munich, University Observatory / National Astronomical Observatories, Chinese Academy of Sciences

In collaboration with: Yang Huang, Zhixuan Kang, Huaxi Chen, Chenzhou Cui, Jashu Pan, Xiao Kong, Xiaoyu Tang, Hengqiang Han, Haiyang Mu, Yuan-sen Ting, Yunfei Xu, Dongwei Fan, Guirong Xue, Ali Luo and Jifeng Liu

federica b. b...  
she/her  
University of Delaware  
Department of Physics and  
Biden School of Public Policy and  
Data Science Institut

NADC National Astronomical Data Center  
国家天文科学数据中心  
2025-06-05 · Athens · UniversAI2025  
之江实验室 ZHEJIANG LAB  
中国科学院 CAS

# Knowledge



## **Syllogism:**

All men are mortal.

Socrates is a man.



Socrates is a mortal.



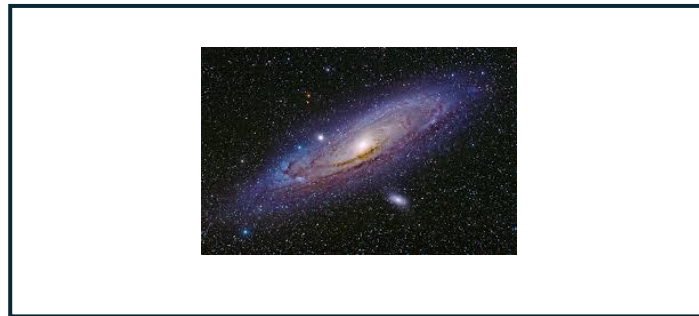
# ... more knowledge



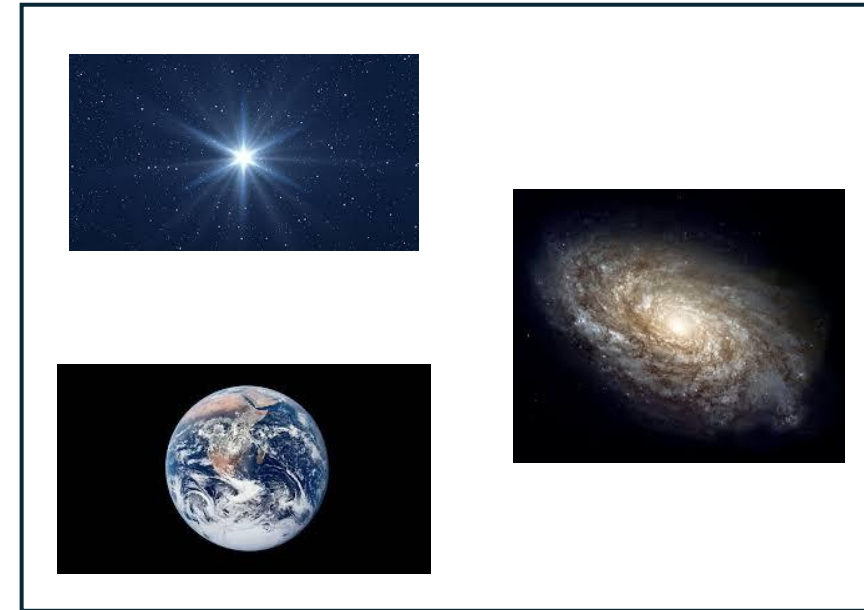
# Knowledge Graphs (KGs)

- **Class:** A group of entities that share common properties or characteristics. i.e., **Celestial bodies**.
- **Subclass:** Specific type of a broader class. i.e., **Galaxies**.
- **Individual/Instance:** Something or someone belonging to a class. i.e., **Milky Way**.

Individual: **Milky Way**



Class: **Celestial Bodies**



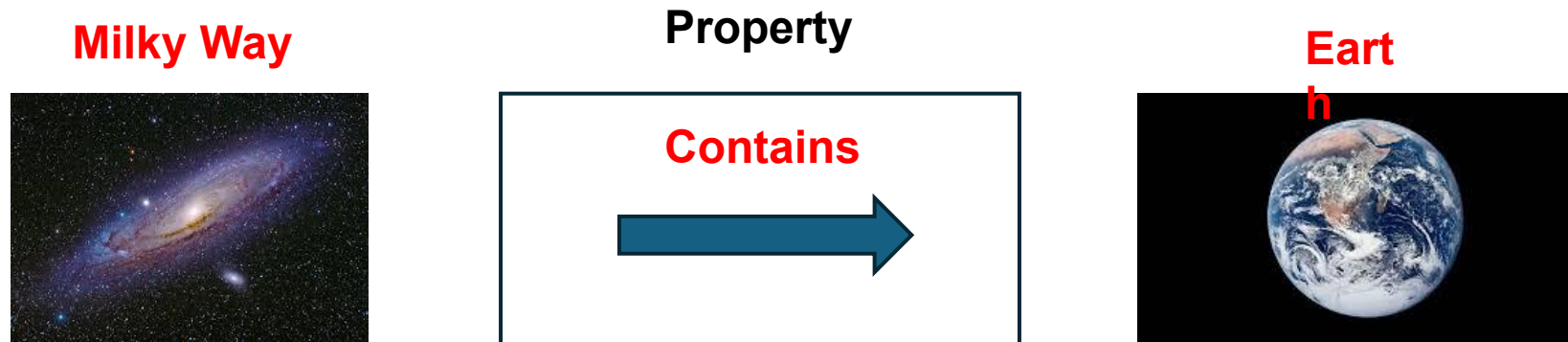
Subclass: **Galaxies**



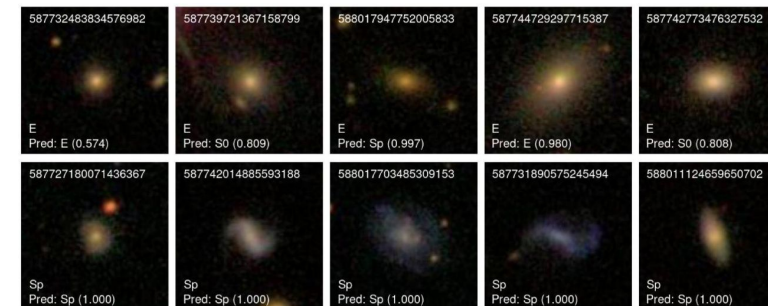
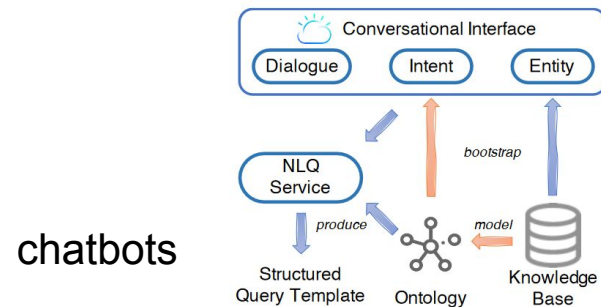
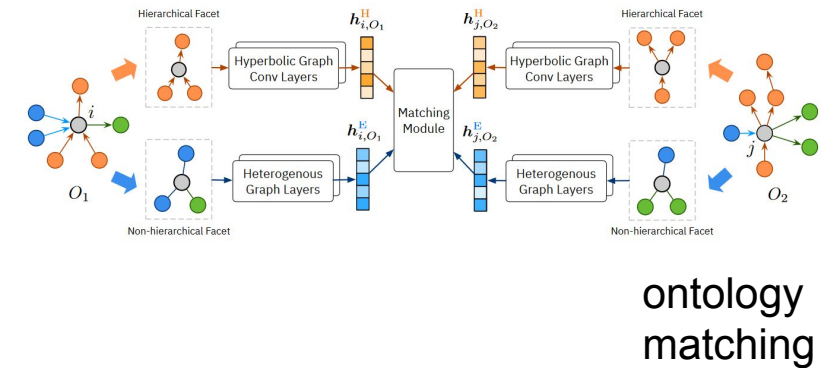
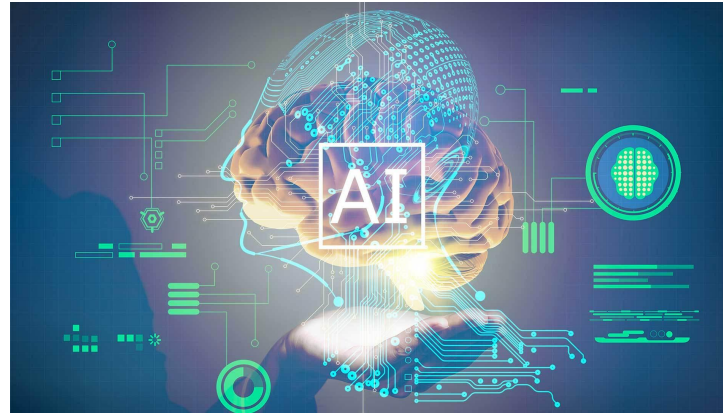
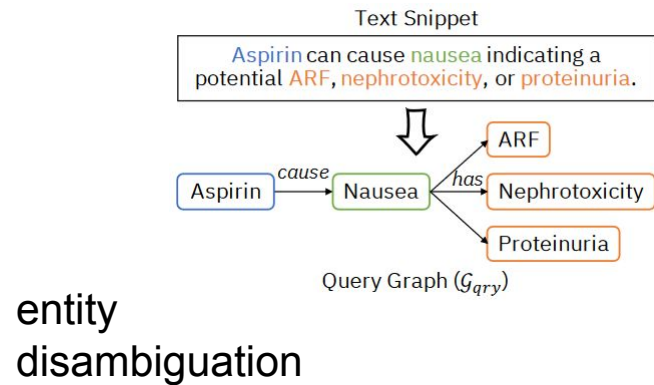


# Knowledge Graphs (KGs) (2)

**Properties:** The attributes that connect individuals.



# Exploiting KGs for deep learning tasks





# Motivation



Table B.1. Spin and external pressure gradient estimates for blazars. Columns:

Table B.1. Spin and external pressure gradient estimates for blazars. Columns:

Table B.1. Spin and external pressure gradient estimates for blazars. Columns:

Table B.1. Spin and external pressure gradient estimates for blazars. Columns:  
(1) Name as given in H09, L17, (2) Alternative name, (3) Class (B is for BL  
Lacs, F for FSRQs), (4) Lorentz factor, (5) Spin, (6) External pressure gradient,  
(7) Reference for the Lorentz factor estimate.

Name	Alt-name	Class	$\Gamma$	$a$	$s$	Ref.
J0003-066	NRAO 5	B	3.3	0.9	0.38	H09
J0016+731	-	F	6.8	0.59	0.64	H09
J0102+5824	0059+5808	F	12.0	0.6	0.83	L17
J0106+013	OC 012	F	27.8	0.82	1.07	H09
J0136+4751	0133+476	F	9.5	0.49	0.76	L17
J0202+149	4C 15.05	F	9.9	0.61	0.76	H09
J0212+735	-	F	7.5	0.57	0.67	H09
J0217+0144	PKS 0215+015	F	19.1	0.75	0.96	L17

read tables

June 18, 2025



cross-check online  
astronomical catalogues and  
DBs

SummerSoC 2025, Hersonissos, Crete, Greece

## 9286 Stars: An Agglomeration of Stellar Polarization Catalogs

Carl Heiles  
Astronomy Department, University of California, Berkeley, CA 94720

### ABSTRACT

This is a revision. The revisions are minor. The new version of the catalog should be used in preference to the old. The most serious error in the older version was that  $\theta_{eff}$  was incorrect, being sometimes far too large, for Reiz and Franco entries; the correct values are all zero for that reference.

We present an agglomeration of stellar polarization catalogs with results for 9286 stars. We have endeavored to eliminate errors, provide accurate ( $\sim$ arcsecond) positions, sensibly weight multiple observations of the same star, and provide reasonable distances. This catalog is available by anonymous FTP as ascii file <ftp://vern1.berkeley.edu/pub/polecat/p14.out>. This manuscript is also available as the postscript file <ftp://vern1.berkeley.edu/pub/polecat/po11.ps>.

Subject headings: catalogs — ISM: magnetic fields — ISM: dust, extinction — stars: distances

### 1. INTRODUCTION

Polarization has been measured for thousands of stars and presented in perhaps a dozen catalogs. Some previous attempts to combine these lists are very admirable because they have made it much easier to use the data. The largest include Mathewson et al (1978; hereafter MFKNK) catalog (CDS catalog I/34A) and Axon and Ellis (1976) (CDS catalog I/178). However, they have deficiencies; for example, both list multiple results for individual stars and have not purged errors from the original catalogs. The present agglomeration combines multiple observations with weighted averages, fixes most errors, provides accurate positions, and reasonable estimates for stellar parameters such as distance and extinction. It also includes information on which original catalogs were used for each entry.

Section 3 discusses the catalogs that we have included, together with the information contained in each. The MFKNK, Axon and Ellis (1976), Reiz and Franco (1998) and Goodman (1997) catalogs were originally provided to us in electronic form. We entered the Appenzeller (1974) catalog by hand from the printed page. For all the other catalogs, we scanned printed

<sup>1</sup>email: [cheiles@astro.berkeley.edu](mailto:cheiles@astro.berkeley.edu)

publish combined data  
and findings on them



combine  
data





HECATE

# The **H**eraklion **E**xtragalactic **CAT**alogue**E**



# The good news...

other query modes :

Identifier query

Coordinate query

Criteria query

Reference query

Basic query

Script submission

TAP

Output options

Help

Query : sirius

submit id

Basic data :

\* alf CMa -- Spectroscopic Binary

Other object types:

ICRS coord. (ep=J2000) :

FK4 coord. (ep=B1950 eq=1950) :

Gal coord. (ep=J2000) :

Proper motions mas/yr :

Radial velocity / Redshift / cz :

Parallaxes (mas):

Spectral type:

Fluxes (8) :

SIMBAD Query around within 2 arcmin

06 45 08.91728 -16 42 58.0171 (Optical) [ 11.70 10.90 90 ] A 2007AB...474..653V

227.23029126 -08.89028121 [ 11.70 10.90 90 ]

-546.01 -1223.07 [1.33 1.24 0] A 2007AB...474..653V

V(km/s) -5.50 [0.4] / z(-) -0.000018 [0.000001] / cz -5.50 [0.40]

A 2006AstL...32..759G

379.21 [1.58] A 2007AB...474..653V

A0mAlVa C 2003AJ...126.2048G

U -1.51 [-] C 2002yCat.2237....00

B -1.46 [-] C 2002yCat.2237....00

V -1.46 [-] C 2002yCat.2237....00

R -1.46 [-] C 2002yCat.2237....00

I -1.43 [-] C 2002yCat.2237....00

J -1.36 [-] C 2002yCat.2237....00

H -1.33 [-] C 2002yCat.2237....00

K -1.35 [-] C 2002yCat.2237....00

All (CDSPortal)

Send to

Photometry within 5 arcsec

## Identifiers (63) :

An access of full data is available using the icon Vizier near the identifier of the catalogue

* alf CMa	GEN# +1.00048915A	LPH 243	ROT 1088
* 9 CMa	GJ 244 A	LTT 2638	SAO 151881
* alf CMa A	HD 48915	2MASS J06450887-1642566	SBC7 288
** AGC 1A	HD 48915A	N30 1470	SBC9 416
ADS 5423 A	HGAM 556	NAME Dog Star	SKV# 11855
AKARI-FIS-V1 J0645085-164258	HIC 32349	NAME Sirius	TD1 8027
BD-16 1591	HIP 32349	NAME Sirius A	TIC 322899250
BD-16 1591A	HR 2491	NLT 16953	TYC 5949-2777-1
CDM J06451-1643A	IDS 06408-1635 A	NSV 17173	UBV 6709

## Sources:

- supplementary material optionally submitted by authors
- Huge manual work of curators
- User-provided feedback

<http://simbad.u-strasbg.fr/Pages/guide/ch02.tx>

# The bad news...

## Existing cross-domain KGs



active galactic nucleus (Q46587)

compact region at the center of a galaxy that has a much higher than normal luminosity over at least some portion – and possibly all – of the electromagnetic spectrum [edit](#)  
Active galaxy | AGN

[In more languages](#)  
[Configure](#)

Language	Label	Description	Also known as
English	active galactic nucleus	compact region at the center of a galaxy that has a much higher than normal luminosity over at least some portion – and possibly all – of the electromagnetic spectrum	Active galaxy

[All entered languages](#)

### Statements


instance of	<a href="#">astronomical object type</a> <a href="#">edit</a>
	<a href="#">0 references</a>
	<a href="#">+ add reference</a>
	<a href="#">+ add value</a>


subclass of	<a href="#">galaxy</a> <a href="#">edit</a>
	<a href="#">0 references</a>
	<a href="#">+ add reference</a>





# The ugly news...


Existing cross-domain KGs:

 DBpedia

 Browse using ▾

 Formats ▾

 Faceted Browser

 Sparql Endpoint

## About: [Active galactic nucleus](#)

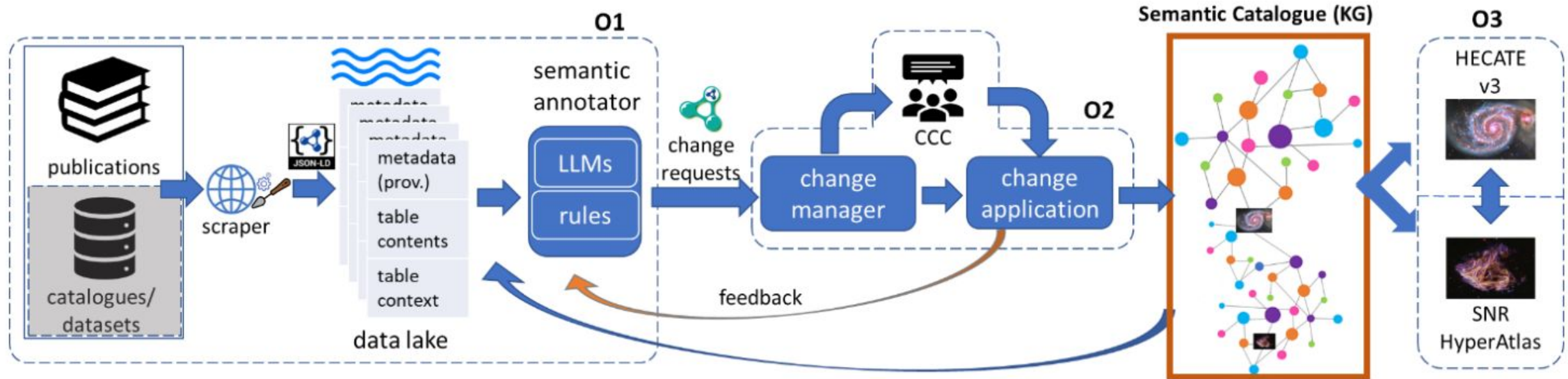
→ An Entity of Type: [settlement](#), from Named Graph: <http://dbpedia.org>, within Data Space: [dbpedia.org](http://dbpedia.org)

An active galactic nucleus (AGN) is a compact region at the center of a galaxy that has a much-higher-than-normal luminosity over at least some portion of the electromagnetic spectrum with characteristics indicating that the luminosity is not produced by stars. Such excess non-stellar emission has been observed in the radio, microwave, infrared, optical, ultra-violet, X-ray and gamma ray wavebands. A galaxy hosting an AGN is called an "active galaxy". The non-stellar radiation from an AGN is theorized to result from the accretion of matter by a supermassive black hole at the center of its host galaxy.





# PARSEC





# DOCBO and HECAT3KG



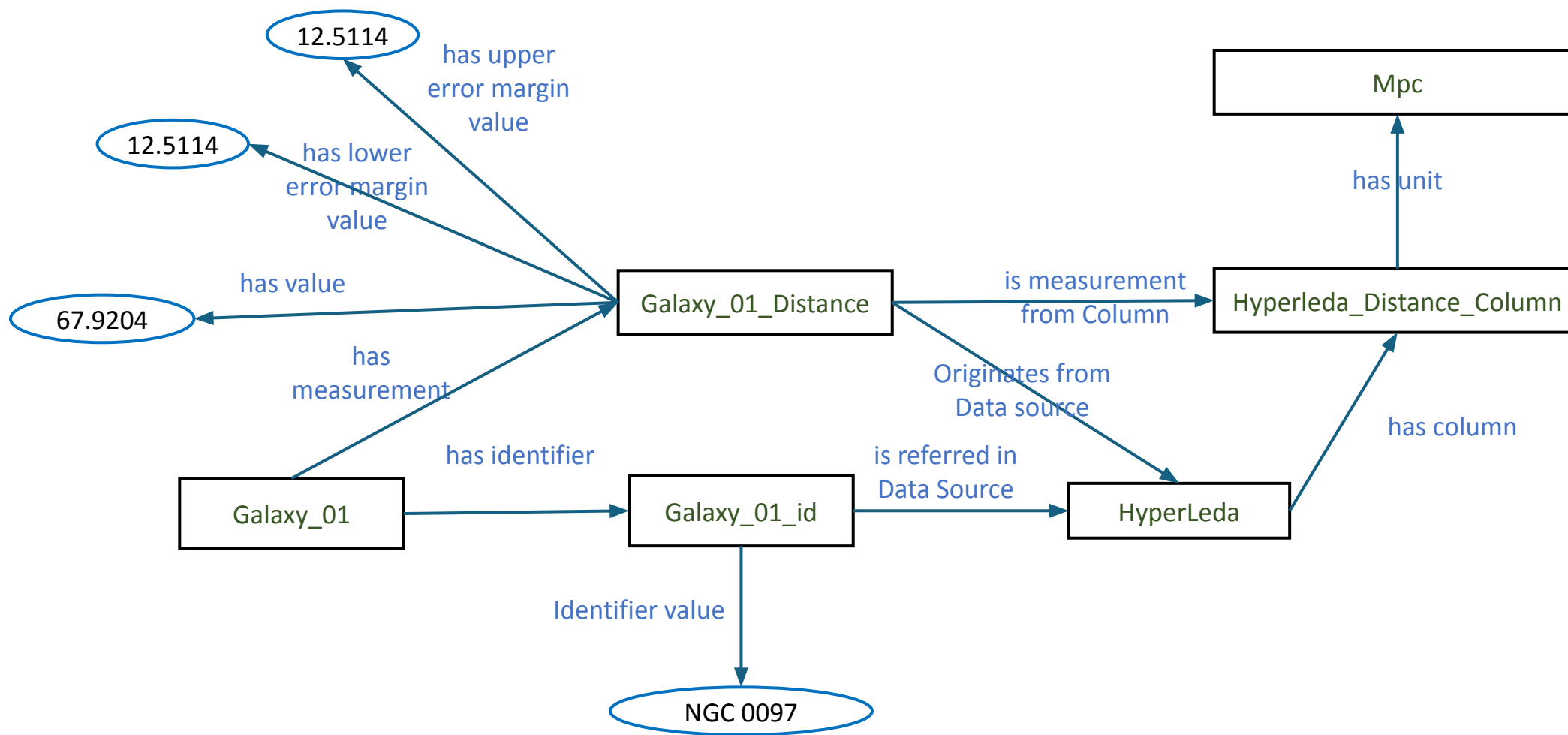
- **DOCBO**: the ontology we created, containing 5 main classes:
  - Celestial Body, Identifier, Data Source, Measurement, Data Column
  - Public: <https://zenodo.org/records/15388573>



- **HECAT3KG**: A KG, following DOCBO, representing HECATE
  - following two mapping approaches (RML and X3ML)
  - 44 million RDF triples (~5GB)
  - Public: <https://zenodo.org/records/15379419> (along with mappings)

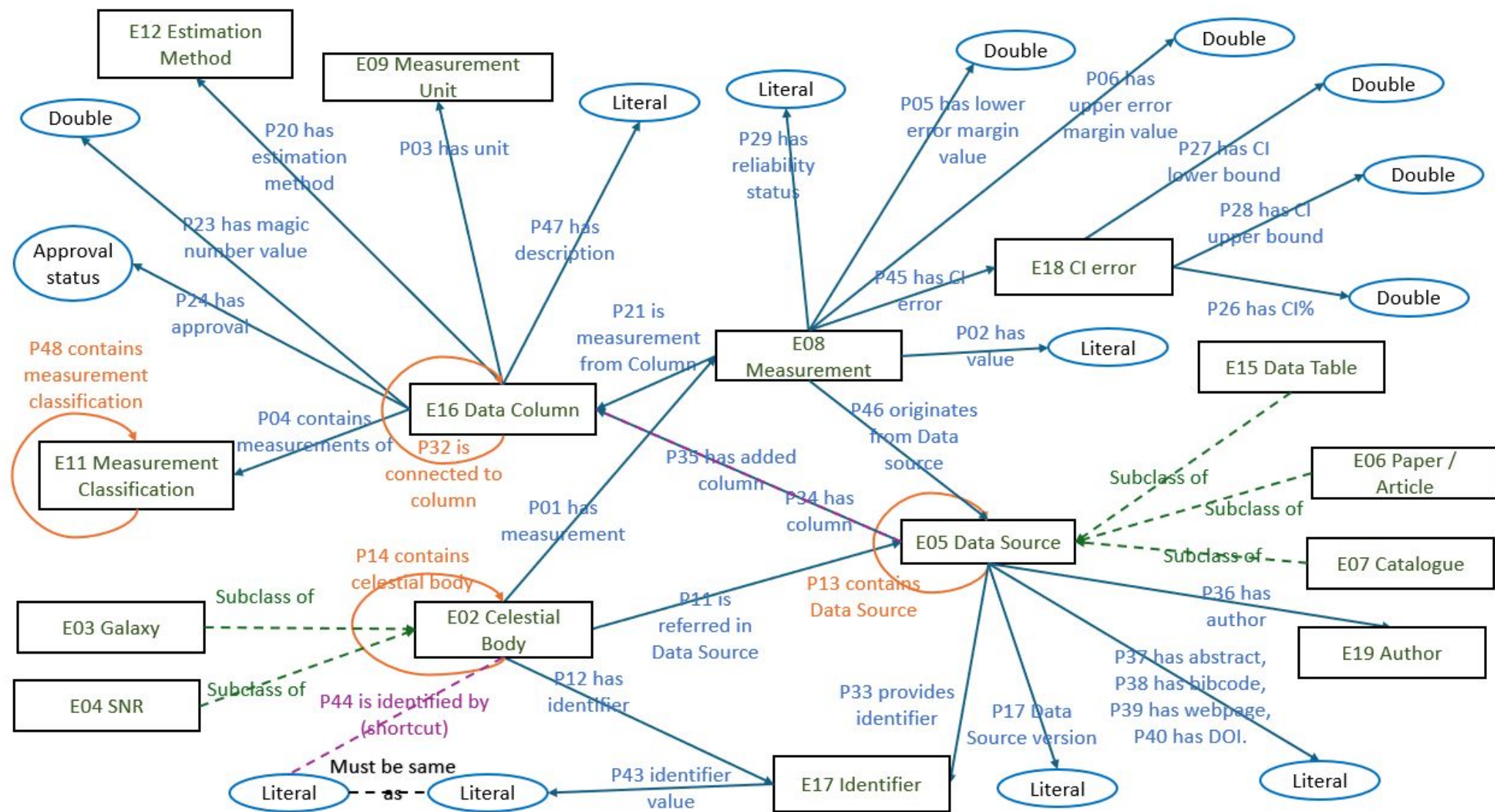
PGC	OBJNAME	ID_NED	...	ID_2MASS	...	RA	DEC	...	D	E_D	...
1442	NGC0097	NGC 0097		00222998+2944433		5.624916	29.745361		67.9204	12.5114	

# HECAT3KG - Example

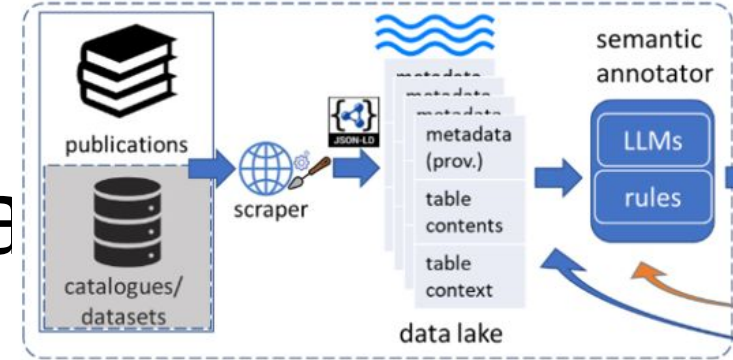




# DOCBO



# Tabular Data Mining and Annotation



- Read pdf/mrt publication files
- Semantic Table Interpretation
  - Table Topic Detection
  - Cell Entity Annotation
  - Column Type Annotation
  - Columns-Property Annotation
- Neuro-symbolic approach
  - Statistical models & LLMs to suggest annotations
  - Rules to reduce the search space

RA: (Pepsi , wd:Q47719)

CEA: (Canada, wd:Q16)

CTA:  
(country wd:Q6256)

E1	Switzerland	Rivella	1952-01-01
E2	U.S.	PepsiCo	1961-01-01
E3	Canada	The Coca-Cola Company	1998-01-01

CPA: inception, wdt:P571

TD: drink, wd:Q40050

# Conclusions

- AI (**Bits**) is no longer just a tool; it's an integral part of astrophysics
- Astrophysics (**Stars**) is entering a golden era of data (e.g., LSST, SKA) that CS can uniquely address
- Ontologies/KGs (**Knowledge**) offer a bridge between raw data and a deep understanding (years of research founded in FOL)
- Interdisciplinary efforts like **UniversAI** show that collaboration isn't just possible—it's urgently needed.



# Thank you!