Challenges and Innovations in Learning from Heterogeneous Data in Fundamental Physics

¹L2I Toulouse, CNRS/IN2P3, Université de Toulouse





C.Biscarat¹, <u>S.Caillou¹</u>, J.Stark¹

EUCAIFCon2025

Université de Toulouse

















AISZAT C. Biscarat, S. Caillou, J. Stark



WorkShop on "Heterogeneous Data and Large Language Model"



C. Biscarat, S. Caillou, J. Stark

~80 AI experts gathered over 4 days in Toulouse, 28 talks including 6 keynote adresses

WorkShop on "Heterogeneous Data and Large Language Model"

Bring people together: Participants from major international physics collaborations (CERN, LISA) and the Toulouse AI ecosystem (3IA ANITI, local industry, HPC stakeholders)

Deliberately interdisciplinary: fundamental physics, Earth observation, biology, fluid mechanics, weather forecasting, cognitive science, and more—encouraging cross-field idea sharing and collaboration

ASSAT C. Biscarat, S. Caillou, J. Stark

~80 AI experts gathered over 4 days in Toulouse, 28 talks including 6 keynote adresses

The **AISSAI center** was created in 2021 as a part of the CNRS National strategy in AI interfaces with Al.

AISSAT C. Biscarat, S. Caillou, J. Stark

Thank you to AISSAI for making this workshop possible!

- Main objective: structure and organize cross-functional actions involving all CNRS institutes at the

Our workshop was a part of the Artificial intelligence for the two infinites AISSAI thematic trimester

Heterogeneous data in particle physics detectors

AISZAT C. Biscarat, S. Caillou, J. Stark

Gravitational Waves sources for LISA

Cosmological sources stochastic backgrounds + bursts or a background generated by cosmic strings ? Amplitude/rate very uncertain.

Stellar-origin-blackhole inspirals (SOBH) O(few) yr⁻¹

Signal will be thousand of Gravitational Waves coming from everywhere from hetereogenous sources (galactic binaries, mergers of massive black holes, ...)

C. Biscarat, S. Caillou, J. Stark

Massive black hole binaries (MBHBs) Rate uncertain, could be several tens per year

Extreme-mass-ratio inspirals (EMRIs) Rate unknown, could be as low as ~1 yr⁻¹ or as high as 1000 yr⁻¹

Ultra-compact binaries (UCBs) ~10⁷ systems ~10⁴ resolvable

Exotic sources

LISA Data Analysis challenge (The Global fit)

Hundreds

Can AI help with that ?

C. Biscarat, S. Caillou, J. Stark

Separating overlapping Gravitational Waves signals is an extremely hard problem (The Global fit challenge)

Heterogeneous Data in Astrophysics

Images, spectra, and time-series measurements of millions of astrophysical phenomena

AISSA[®] L2T C. Biscarat, S. Caillou, J. Stark

The Multimodal Universe: 100TB of **Astronomical Scientific Data**

Goal: Assemble the first large-scale multi-modal dataset for machine learning in astrophysics. Main pillars:

- Engage with a broad community of AI+Astro experts. 0
- Adopt standardized conventions for storing and 0 accessing data and metadata through mainstream tools (e.g. Hugging Face Datasets).
- Target large astronomical surveys, varied types of 0 instruments, many different astrophysics sub-fields.

Polymathic

Can we build AI models which leverage information from heterogeneous datasets and learn better representations of the physic objects?

AISSAT L2T C. Biscarat, S. Caillou, J. Stark

Multimodality Architectures: State of the Art

(Self)

Massively supervised learning (« brute-force »):

- **Representation alignment:** 0
- Image-to-Text (captioning):
- **Text-to-Image:** 0

CLIP (OpenAl)	400M imag
CoCa (Google)	4B
DALL-E3 (OpenAl)	>5B
lmagen (Google)	

AISZAT C. Biscarat, S. Caillou, J. Stark

ge-text pairs

The rise of foundation models

Foundation model approach

- **Combine** pretrained modules in more complex systems.

Foundation model take advantage of **multimodality**:

- Self-supervised training by contrastive cross modal tasks
- Enrich data representation in latent spaces by semantic grounding between modalities

• Pretrain models on pretext tasks, with self-supervision, on very large scale datasets. Adapt pretrained models to downstream tasks (Transfert Learning).

On the Opportunities and Risks of Foundation Models, Bommasani et al 2021

Learning from heterogeneous data in ATLAS

Graph Neural Networks for track reconstruction in the ATLAS ITk detector, Minh Tuan Pham o.b.o GNN4ITk team

Heterogeneous, Multi-Task Models for Flavour Tagging in ATLAS Jackson Barr o.b.o the SALT team

AISSA

GNN deals with heterogeneous inputs from the Pixel and the Strip sub-systems in the ATLAS ITk

Current generation of flavour tagging model in production in ATLAS,

AISSA®

Keynote Adress: Foundation models for high energy physics, Anna Hallin

Biscarat, S. Caillou, J. Stark C.

Focus on OmniJet- α

- Uses a transformer for generative pretraining based on the GPT-1 architecture with next-token-prediction as training target. $p(x_j|x_j-1, \ldots, x_1, < \text{start token} >)$
- Generation : Generally good agreement to truth distribution
- Able to task-switch: unsupervised full jet generation to supervised classification

Jet type prediction

$\mathsf{Jet} = \{p_1, p_2, \dots, p_N\}$

Jets as **sequences of integers**:

{< start token >, token₁, token₂, ... , token_N, ·

AISZAT C. Biscarat, S. Caillou, J. Stark

			OmniJ					
or	HE	Anomaly detection		UCAIF	Con	2025	Cross-entropy class labels + diffusion velocity parameter	
			Highlight talk	Second Se	Models	Highlight talks	Jet classification	
~	de	coder .	$egin{array}{cccc} & & & & & & & & & & & & & & & & & $	(\vec{p}_{i}) x $\vec{p}_{i} = (p_{T}$	$egin{array}{cccccccccccccccccccccccccccccccccccc$	encoder	token n	

Foundation models in Neutrino Physics

AISSAT

Small thinks big: transfer learning in KM3NeT/ORCA for neutrino event reconstruction, Iván Mozún Mateo, Antonin Vacheret o.b.o the KM3NeT collaboration

DIKU, University of Copenhagen

Biscarat, S. Caillou, J. Stark

DL classification Track-like vs Shower-like

> Track and shower reconstruction

Foundation model in KM3NeT?

- learns as the detectors grow
- can handle multiple geometries
- can handle both KM3NeT/ORCA ۲ and KM3NeT/ARCA
- classification \leftrightarrow reconstruction \bullet

AISZAT C. Biscarat, S. Caillou, J. Stark

Keynote Adress: <u>Deep Learning & the Global Workspace Theory</u>, Ruffin Van Rullen

Rufin Van Rullen CNRS, The Brain and **Cognition Research Center** (CerCo), Artificial and Natural Intelligence Toulouse Institute (ANITI), Université de Toulouse, France.

Multimodal architectures: Limits

- - (a child probably gets <<1M of explicit supervision examples)
- Sub-optimal multimodal grounding
- Sub-optimal compositionality

a blue cube on top of a red cube, beside a smaller yellow sphere

- → Solution A: bigger models, trained with even more data?
- → Solution B: change of paradigm?

C. Biscarat, S. Caillou, J. Stark

Devillers et al (CoNLL 2021): Does language help generalization in vision models?

The Global Workspace Theory

hierarchy of modular processors

Bernard J. Baars (1993) former Senior Fellow in Theoretical Neurobiology at the Neurosciences Institute in San Diego, US.

🟦 Access through **your organization**

Purchase PDF

Progress in Brain Research Volume 150, 2005, Pages 45-53

automatically activated processors

Global workspace theory of consciousness: toward a cognitive neuroscience of human experience

Bernard J. Baars 😤 🖾

Show more \checkmark

+ Add to Mendeley 😪 Share 🍠 Cite

https://doi.org/10.1016/S0079-6123(05)50004-9 7

Get rights and content 7

AISZAT C. Biscarat, S. Caillou, J. Stark Al for science, science for Al

Deep learning and the Global Workspace Theory

AISZAT C. Biscarat, S. Caillou, J. Stark

🟦 Access through your organizatior

View Open Manuscript

Purchase PDF

Trends in Neurosciences

Volume 44, Issue 9, September 2021, Pages 692-704

Deep learning and the Global Workspace Theory

Rufin VanRullen ¹² $\stackrel{\frown}{\sim}$ 🖾 , Ryota Kanai ³

- ¹ The Brain and Cognition Research Center (CerCo), CNRS UMR5549, Toulouse, France
- ² Artificial and Natural Intelligence Toulouse Institute (ANITI), Université de Toulouse, France
- ³ Araya Inc., Tokyo, Japar

Towards implementating GWT

Modules

(hundreds of available choices!)

Output Choice determines model functionality

- Our Consumption Of Consumpti Consumption Of Consumption Of Consumption Of Cons
- Trained via cycle-consistency objective

Output Description

- Transformer: key-query matching
- Top-down & bottom-up control

Multimodal systems with Global Workspace

Our ecosystem...

Funded by the European Union

European Research Counc

ERC Advanced Project GLOW (2023-2028)

- Develop brain-inspired multimodal deep learning systems
- Evaluate their use and relevance for machine learning
- Advance our knowledge of the brain

ANITI Synergy Chair C3-PO (2024-2028)

- Cobots with Conversation, Cognition & Perception
- Chairs: R.VanRullen (CerCo), N. Asher (IRIT), T. Serre (Brown), O. Stasse (LAAS)
- Frugal multimodal robotic systems with grounded perception, language and action

Léopold Maytié

Ph.D Student CerCo Université de Toulouse Toulouse (France)

in M

 $\lambda_1.\lambda_2 = \text{softmax}(K_1.0, K_2.0)$

LISA Data Analysis challenge: Towards an Al-driven Global fit?

- Successful solutions to the global fit problem have used classic stochastic sampling techniques.
- Typical strategy adopted is to iteratively update the solution for one source type and then move to the next.

Scaling and computing ressources stay a challenge

softmax(QK

keytask

key_{la}

Simulation Based Inference (SBI) is explored to learn posterior for a single source of data

<u>Neural density estimation for Galactic</u> Binaries in LISA data analysis Natalia Korsakova (2024)

Would it be possible to combine type of sources dedicated neural density estimators in a "Global Workspace" with Attention Mechanisms ??

Original idea from J. Stark

Why choose a hierarchical, modular architecture (à la Global Workspace) over a single giant, fully-tokenized Transformer?

Preserves symmetry & equivariance

Each modality keeps its own specialist encoder that respects the intrinsic geometric / group-theoretic structure of the data (e.g., CNNs for images, SE(3) GNNs for molecules), instead of flattening everything into one generic token space.

Built-in explainability

Clear boundaries between modules + a central workspace let you trace which specialist contributed what evidence and when, making reasoning chains inspectable rather than buried in millions of mixed-modality parameters.

Future-proof modularity

New capabilities—Retrieval-Augmented Generation, neuro-symbolic reasoning, causal-learning engines—can be plugged in as standalone processors that broadcast to / read from the workspace, without retraining the entire stack.

extensibility better than an ever-larger monolith.

A structured "society of experts" scales knowledge, transparency, and

Thanks to all the keynotes !

Anna Hallin

Institute of Experimental Physics, Universität Hamburg (Germany). Particle physicist. Lead of the effort on Foundation Models in the research group for machine learning in particle physics at Hamburg. \rightarrow biography

Jordi Inglada

French Space Agency (CNES, Toulouse), Center for Spatial **Biosphere Studies (CESBIO,** Toulouse). Main researcher in the "Artificial and Natural Intelligence Toulouse Institute" (ANITI) in the theme "Learning with little complex data, AI and phys models". \rightarrow biography

François Lanusse CNRS, UMR AIM / Flatiron Institut.

Michel Besserve

Empirical Inference Department, Max Planck Institute for Intelligent Systems, Tübingen, Germany. Centers of interest: Causality and Machine Learning For Complex Systems. \rightarrow biography

AISZAT C. Biscarat, S. Caillou, J. Stark

Andrew El-Kadi

Research Engineer at Google DeepMind working as part of the Weather effort. \rightarrow biography

Jonathan Gair

Max Planck Institute for **Gravitational Physics - Albert** Einstein Institute - (AEI Potsdam, Germany). Group leader in Gravitational Waves data analysis (current LIGO and PTA detectors, future spatial mission LISA and ground-based Einstein Telescope apparatus).

Rufin VanRullen

CNRS, The Brain and Cognition Research Center (CerCo), Artificial and Natural Intelligence Toulouse Institute (ANITI), Université de Toulouse, France.

nature

Publish with us About the journal \sim Explore content ~

nature > articles > article

Article Open access Published: 04 December 2024

Probabilistic weather forecasting with machine learning

<u>Ilan Price</u> ⊠, <u>Alvaro Sanchez-Gonzalez</u>, <u>Ferran Alet, Tom R. Andersson, Andrew El-Kadi, Dominie</u> Masters, Timo Ewalds, Jacklynn Stott, Shakir Mohamed, Peter Battaglia 2, Remi Lam 2 & Matthew

Nature 637, 84–90 (2025) Cite this article

Thanks to all the participants and the committees!

Scientific advisory committee

- Sylvain Caillou (L2IT, IN2P3, CNRS/UT3) chair
- Alexandre Boucaud (APC, IN2P3, CNRS)
- Tobias Golling (Université de Genève)
- François Lanusse (Polymathic AI)
- Daniel Murnane (Copenhagen University)
- Thomas Oberlin (ISAE-SUPAERO, ANITI, Université de Toulouse)
- Jan Stark (L2IT, IN2P3, CNRS/UT3)
- Gordon Watts (Washington University)

Local organisation committee

- Catherine Biscarat (L2IT, IN2P3, CNRS/UT3) chair
- Sylvain Caillou (L2IT, IN2P3, CNRS/UT3)
- Jocelyne Gauthier (L2IT, IN2P3, CNRS/UT3)
- Jan Stark (L2IT, IN2P3, CNRS/UT3)
- Jeanette Thibaut (L2IT, IN2P3, CNRS/UT3)

Scientific secretaries

- Vasco Gennari, PhD student (L2IT, CNES, CNRS/IN2P3, UT3)
- Alexandro Martone, PhD student (L2IT, IN2P3, CNRS/UT3)
- Minh-Tuan Pham, PhD student (University of Wisconsin-Madison)

catherine.biscarat@l2it.in2p3.fr

Do not hesitate to get in touch with us !

jan.stark@l2it.in2p3.fr

sylvain.caillou@l2it.in2p3.fr

EuCAIFCon2025

