## FOUNDATION MODELS AND AGENTS FOR PHYSICS

Nesar Ramachandra nramachandra@anl.gov CPS Division, Argonne National Laboratory





## INTRODUCTION



- What are Foundation Models (FMs)? What makes them particularly interesting or powerful?
- What does it take to train a Large Language Model (LLM)? How can we align the LLMs to our science domains?
- How much productivity can we harness out of LLMs? Can we go beyond text-type data? Can we use a system of AI-agents/experts in solving physics problems?
- All the applications here are outside of nuclear physics, however they are adaptable.

\*Disclaimer: Some of the opinions are my own.



## **HOW ARE FOUNDATION MODELS DIFFERENT?**



### **Traditional AI:**

- Parameters: O(1e2-1e8)
- Data: Specialized, annotated
- *Flexibility:* Meant for specific tasks (classification, regression, generative, unmasking, ...)
- Training Approach: Supervised, unsupervised, or transfer learning.
- *Applications:* Suited for narrowly defined problems with clear objectives.

### Foundation models:

- Parameters: O(1e8 to 1e11)
- Data: Extensive, datasets from various domains
- Flexibility: Highly versatile, can be fine-tuned different tasks.
- Training Approach: Generic pre-training followed by alignments
- Applications: Suitable for multi-faceted problems requiring nuanced understanding



Image credits: Emil Constantinescu & ChatGPT



## **DO FOUNDATION MODELS GENERALIZE?**

- Majority of the foundation models are language based.
- Emergent behaviors seen with larger models.
- Questions about whether this is the path towards AGI, whether there is a reasoning/thinking involved at inference.
- Regardless, FMs are technological paradigm shifts, and a highly useful knowledge base.



Argor

#### Latent space of LLMs showing emergent behavior



## **STAGES IN TRAINING AN LLM**

- Data preparation
  - · Both labeled and unlabeled datasets.
- Pre-training
  - Computationally expensive (10-100 TB data, 10+ million GPU hours, \$100 million)
- Tuning
  - Smaller task, includes fine-tuning or instruction tuning often with supervised datasets.
- Reinforcement learning
  - Safety guard-rails
  - Human-in-the-loop
- Inference improvements
  - Can include reinforcement learning, Chain-ofthought optimization, reasoning tokens and feedback preference optimizations.
  - Heavily utilized in "Reasoning" class models (Open Al o1, o3, ...).
- Benchmarking and validation
- Deployment and API

Sequence of tokens:  $U = \{u_1, u_2, u_3, \dots, u_N\}$ 

Minimize: 
$$L(U) = -\sum_{i=1}^{N} log(p(u_i | u_{i-k}, ..., u_{i-1}, \Theta))$$

(Conditional probability of i-th token given k preceding tokens)

Transformer-based architecture and attention are crucial for scaling training to billions of tokens



#### LLM construction stages: pre-training to reasoning abilities







## **ALIGNING MODELS TO SCIENCE**

#### • Fine-tuning

- Smaller, specific scientific datasets (e.g: AstroLLaMa)
- Full model building
  - Expensive pre-training from scratch (e.g: AuroraGPT at Argonne)
    - Scientific datasets
    - Scientifically rigorous benchmarks

**Completed by LLaMA-2** 

The Magellanic Stream (MS) - an enormous ribbon of gas

spanning 140° of the southern sky trailing the Magellanic

Clouds - has been exquisitely mapped in the five decades

since its discovery. However, despite concerted efforts, no

identified. This stellar stream would reveal the distance and

6D kinematics of the MS, constraining its formation and

stream in the MS, with a 6D velocity of 100 km s-1, and a

distance of 100 kpc. The stream is 100 pc wide and 1000

pc long, and is consistent with the MS. The stream is 1000

times more massive than the Magellanic Clouds, and is

likely to be a tidal stream from the Large Magellanic

Cloud.

the past orbital history of the Clouds. We find a stellar

stellar counterpart to the MS has been conclusively

- Continual pre-training
  - Middle ground.
  - Less annotated but more domainspecific data



#### Online Communities Reference Autos & Vehicles Arts & Entertainment News Sports Visual Art & Design **General Reference** Politics Sports News Local News Business & Industria Books & Literature Internet & Telecon Science Design lumanities Humor liological Scien Web Service: History Online Media People & Society Online Image Galleries Sensitive Subjects Jobs & Education Travel Adult **Computers & Electronics** Porn Education rogramming Software Games Computer & Video Games Law & Government Shopping Health https://arxiv.org/pdf/2204.02311

#### Data-corpus division in LLMs

#### Fine-tuning model performance

#### Completed by AstroLLaMA

The Magellanic Stream (MS) - an enormous ribbon of gas spanning 140° of the southern sky trailing the Magellanic Clouds - has been exquisitely mapped in the five decades since its discovery. However, despite concerted efforts, no stellar counterpart to the MS has been conclusively identified. This stellar stream would reveal the distance and 6D kinematics of the MS, constraining its formation and the past orbital history of the Clouds. We 1) present a new method to select MS candidate stars from wide field optical imaging, 2) apply this method to the MS, 3) compare the results to existing catalogs, and 4) identify the most promising MS candidate stars. We apply our method to 2547 stars in the Gaia EDR3 catalog that have Gaia proper motions consistent with the MS and use these as a pilot sample to optimize our method. We then apply our method to the entire Gaia EDR3 catalog to identify 46,000 stars that are candidates for the MS, and we use the Gaia-ESO Survey to identify 13,000 of these stars that are likely to be MS candidates. We identify 266 stars with (G-K) colours and [Fe/H] that are consistent with previous estimates of the MS. These stars have proper motions and radial velocities that are consistent with the MS, and we place them on a 3D map of the MS. We find a previously unknown stellar stream, which we call the Southeastern Stream (SES), in the direction of the Clouds and MS. The SES has an M31 distance of 32.2+/-0.5 kpc and a 6D velocity of 136+/-10 km/s. The SES is about 20 times brighter in stars than the SMC. The SES is also 2.7 times more metal-rich than the MS which suggests that it is a separate population.

## **LLM-FOR ASTROPHYSICS: ASTROSAGE**

- Continual pre-training + Finetuning
  - Texts from astrophysics papers as training tokens. (20 GB data)
  - 2nd step of fine-tuning the models to answer astrophysics questions (10 GB data)
  - Built on open-source Llama-8B.
  - For domain-specific tasks, as powerful as GPT-4o, with fraction of the training cost (30k GPU hours).



Model cost vs performance: AstroSage vs rest.



## **BEYOND-TEXT: FULL MODEL BUILDING**

- Foundation models for non-text scientific data formats:
  - For science datasets: <u>Star-foundation models</u>, <u>AstroCLIP</u>, <u>OmniLearn</u>
- Numbers be processed as numbers and not as tokenized representations: <u>xVal</u>
- Multi-modality in object expressions, not just data-types.
  - Distinct from industry trend: videos, subtitles, plots and captions, songs and lyrics
  - Nuclear physics: particle trajectories (graphs), events (scalars), catalog values, metadata, instrument settings.
- Variety of transformers
  - Vision Transformers, graph transformers
- Latent space modeling using Contrastive learning

Multiple modalities in astrophysics data. Galaxy distribution in the sky (left), randomly chosen spectra (bottom left), images (right), catalog values (not shown)









## **DOMAIN FOUNDATION MODELS**



- Training different (encoder-decoder) transformers with different data modalities. Numerical tokenization used.
- Contrastive loss to connect simulation modes with observed/experimental modes
  - Changes embedding space



11.5 12.0 12.5 13.0 13.5 log<sub>10</sub>(*M<sub>H</sub>M*<sup>-1</sup><sub>\*</sub>) (SFH)

correlated

11.5 12.0 12.5 13.0 13.5  $\log_{10}(M_H M_*^{-1}) \text{ (mag)}$ 

## **DOMAIN FOUNDATION MODELS**



- Training is general purpose, deployment is task-specific.
- Flexibility in deployment: queries dictate latent space access.
- Compatibility wrto datasets in multiple domains
- DFMs can be joined with existing LLMs for contexts along with knowledge base access



Flexible input/output combinations: Single input predicting multiple outputs.



Flexible input/output combinations: Multiple inputs used to predict the same output

DOI: 10.22541/essoar.172675995.55091022/v1

## **MULTI-AGENT FRAMEWORKS**



Ingredients: tools, servers, and databases

- LLMs or Reasoning models with tool access (web search, data base reading)
- Multiple expert interactions (Roles: Manager, domain-expert, coding expert, critic)
- Local datasets (texts and numerical data) access
  - Recent datasets, codes, publications
  - Proprietary data
- Local simulation interface access.
- Python server to run codes, feed results back to LLMs



#### Agentic framework of Dr. MACS (Multi-Agent Collaborative System)

## **DR MACS IN ACTION: SCIENCE QUESTION**





New run

What is the average brightness of galaxies with stellar mass between 10^10 and 10^12 solar masses and within the redshift range 0.1<z<0.2?

HfApiModel – Qwen/Qwen2.5-Coder-32B-Instruct –

User inputs a question

## **DR MACS IN ACTION: MAIN AGENT**



Step 1	
Output message of the LLM:	
Thought: To solve this task, I need to retrieve data for galaxies within the specifie an astronomical catalog. Then, I will calculate the average brightness of these galax retrieval from astronomical catalogs, so I will use it for this task.	d stellar mass and redshift range from ies. The CodingAgent can help with data
CodingAgent('query string: Retrieve data for galaxies with stellar mass between 10^10 redshift range 0. 1 <z<0. 2="" catalog')<br="" from="" galaxy="" sdss="" the="">```<end_code></end_code></z<0.>	and 10^12 solar masses and within the
	Main Agent plans the tasks
New run	
You're a helpful agent named 'CodingAgent'. You have been submitted this task by your manager. Task: query string: Retrieve data for galaxies with stellar mass between 10^10 and 10^12 solar range 0. 1 <z<0. 2="" catalog<br="" from="" galaxy="" sdss="" the=""> You're helping your manager solve a wider task: so do not just provide a one-line answer information as possible to give them a clear understanding of the answer Your final_answer WILL HAVE to contain these parts: ### 1. Task outcome (short version): ### 2. Task outcome (extremely detailed version):</z<0.>	in Agent delegates the s, provides instructions ar masses and within the redshift er, instead give as much
Put all these in your final_answer tool, everything that you do not pass as an argument And even if your task resolution is not successful, please return as much context as po act upon this feedback.	t to final_answer will be lost. ossible, so that your manager can

## **DR MACS IN ACTION: SECONDARY AGENT**





## **DR MACS IN ACTION: DATA DOWNLOAD**



Execution logs:											1	
objID	ra	dec	redshift	expRad_r		run2d	ellipticity	distance	M_r			
stellar_mass												
0 1237648720141746438	145.186672	-0.893394	0.142611	2.528442	•••	26	0.533659	610.768460	-21.283233			
1.214124e+11	145 100070	0.000004	0 1 105 67	2 520442		26	522650	C10 5701C0	21 202560			
1 123/648/20141/46438	145.1866/2	-0.893394	0.142567	2.528442	•••	26	0.533659	610.579162	-21.282560			
1.2133720+11	150 170351	_0 8780/0	0 136/3/	2 117338		26	0 317376	58/ 310010	-23 0717/0			
6-304723e+11	130.470334	-0.070049	0.120424	2:44/550	•••	20	0.51/520	004.010919	-23:071740			
3 1237648720144040175	150.383283	-0.939088	0.135132	2.838660		26	0.372125	578,736492	-21.052064			
9.812863e+10												
4 1237648720144040025	150.470354	-0.878049	0.136351	2.447338		26	0.317326	583.959305	-23.070433			
6.297137e+11										First sat	of data down	hahcol
					•••					1131 301		loaueu
	171 100000	0 00 4005	0 100500	2 200074		26	0 000700	420 444420	24 740645	from	online databa	ase
95 123/648/2122/022436	1/1.492665	-0.094285	0.102538	2.0968/1	•••	26	0.223/99	439.144130	-21./18015			
1.0130/70+11	158 213610	_0 011631	0 115113	3 056103		26	0 377100	102 000275	_21 1/1671			
1_065709e+11	130.213010	-0.911031	0.113113	2.020102	•••	20	0.377400	492.999273	-21.1410/1			
97 1237648720679927941	148,133192	-0.524590	0.145860	2.376235		26	0.160101	624,681399	-22.819223			
4.996422e+11												
98 1237648720156557509	179.071941	-0.975031	0.109721	2.891736		26	0.747597	469.906262	-20.620606			
6.594965e+10												
99 1237648720156557509	179.071941	-0.975031	0.109727	2.891736		26	0.747597	469.934528	-20.620737			
6.595758e+10												
[100 FOWS X 22 COLUMNS]												

## **DR MACS IN ACTION: IDENTIFYING ISSUES**



AND s.z > 0.1 AND s.z < 0.2 AND p.dered i < 18 AND p.expRad r > 2.0 AND p.ra IS NOT NULL AND p.dec IS NOT NULL AND s .z IS NOT NULL AND p.expRad r IS NOT NULL AND p.dered i IS NOT NULL Code execution failed at line 'print(filtered\_galaxies[['id', 'ra', 'dec', 'redshift', 'stellar\_mass']])' due to: KeyError:"['id'] not in index" [Step 1: Duration 23.71 seconds] Input tokens: 6,566 | Output tokens: 323] Step 3 Output message of the LLM: ndicating that the expected columns are not in the output, it appears the `search\_sdss\_galaxies` Thought, Ci tool might not return all the columns we assumed. We need to refine how we retrieve and filter the data. I will first retrieve the date without any specific columns and then filter it based on redshift and stellar mass. If the stellar mass data is not available directly, I'll need to filter only on redshift and provide a detailed report on the available data and whether we can get the required information. Let's fetch the data again without specifying the columns and analyze it. Code: ``py galaxies = search\_sdss\_galaxies(min\_redshift=0.1, max\_redshift=0.2, num\_entries=100) print(galaxies) `<end\_code> Observation: \*\*Agent wants to execute the following code:\*\* ```python galaxies = search\_sdss\_galaxies(min\_redshift=0.1, max\_redshift=0.2, num\_entries=100) print(galaxies) Press Enter to execute, or type 'no' to abort:

Plain tool calling not providing desired output, Coding agent figuring out a workaround

## **DR MACS IN ACTION: DEBUG AND RE-RUN**





## **DR MACS IN ACTION: FINAL STEPS**





## **DR MACS IN ACTION: FULL RUN**

Multi-Agent Collabora	ative System for Astr	ronomy & Cosmology	
	Dr. MACS		
SME A	gent Code Agent	Image Gen	
 DuckDuc Tool	ckGo SDSS Search	Text2Img Tool	
Visit V Page To	Web GSWLC Search	Download ImgsTool Plot Gals Tool	

## **CONCLUSIONS AND FUTURE OUTLOOK**



- The foundation models have finally facilitated a truly large and deep representations of some of the largest datasets.
- The FMs are mostly trained for language tasks. Encoding non-text information (either directly or indirectly) is the next important step for scientific impact.
- Agentic frameworks allow for automated explorations of data based on hypothesis and planning, code deployment and feedback loops — all geared towards solving a single problem.
- Exciting times ahead, keeping up with the progress is key.

## Questions?



# BACKUP: BUILDING BLOCKS OF FOUNDATION

- Transformer blocks are the fundamental blocks, selfattention is essential functionality:
  - From "Attention is all you need: <u>https://</u> arxiv.org/abs/1706.03762"
- Successor to Recurrent Neural Networks (RNNs), designed to predict next occurrence in a time series.
- Self attention is key to process large amount of data at once.
- Associated 'technical' details: Tokens, Token embeddings, Position encoding, Vocabulary, Context, Context window lengths, Multi-head attention.

Office of

Science

• The reasoning-class models (OpenAl O-models and others) also use reinforcement learning in strategizing steps using reward mechanisms.





(Conditional probability of i-th token given k preceding tokens)

## BACKUP: PHYSICS BENCHMARKS FOR LANGUAGE MODELS

- Multiple choice questions and Free-form questions can be used to evaluate and compare models
- Compared to biology/chemistry, available LLM benchmarks for physics/astronomy is inadequate
  - Both in terms of quantity and quality

Where do most short-period comets come from and how do we know?

Say the pupil of your eye has a diameter of 5 mm and you have a...

Why isn't there a planet where the asteroid belt is located?

Why is Mars red? mmlu-astronomy

Electromagnetic radiation provides a means to probe aspects of the physical Which of the following statements regarding radiation spectra is NOT correct mmlu-college_physics-original-neg	universe. ?	[ "Lines in the infrared, visible, and ultraviolet regions of the s primarily the nuclear structure of the sample.", "The wavelengths i absorption spectrum of an element are among those in its emission s "Absorption spectra can be used to determine which elements are pre stars.", "Spectral analysis can be used to identify the composition dust." ]	pectrum reveal dentified in an pectrum.", sent in distant of galactic
One end of a horizontal, massless spring is attached to a wall. A mass of 0. attached to the other end of the spring and rests on a table. The mass is di	30 kg is splaced…	[ "0.24 mJ", "0.38 mJ", "0.48 mJ", "0.75 mJ" ]	
Things that are equivalent according to the equivalence principle are mmlu-conceptual physics-dev	[ "space twin.",	e and time.", "a traveling twin and a stay-at-home "gravity and acceleration.", "mass and energy." ]	
Colors in a soap bubble result from light	[ "convo "interfo	erted to a different frequency", "deflection", erence", "polarization" ]	





## **BACKUP: BENCHMARKING FOR AURORA-GPT**

- Benchmark development may be crucial for a Science-focussed GPT comparison with existing LLMs.
- Benchmarking team at the Aurora-GPT collaboration has released a web-form to collect science questions of interest — with real-time evaluation from multiple LLMs.
- Goal is to collect O(1000) questions across scientific fields — HELP needed! (and potential collaboration opportunities)
  - High-quality.
  - Should represent what the science community wants out of an LLM.
  - Should not be exposed to current LLMs.

<u>https://web.cels.anl.gov/</u> projects/auroragptquestions/ <u>ui/index.html</u>

## Trial Q:For the force between quarks, which of the following statements is true?

a) The force follows an inverse square law
b) Approximately an inverse square law but asymptomatically weaker at short distances
c) Approximately an inverse square law but asymptomatically stronger at short distances
d) Approximately an Inverse linear relationship, but goes asymptomatically free at large distances.
e) Asymptomatically free at large distances and small distances, inverse relationship in between.



## **BACKUP: SEMI-AUTOMATED BENCHMARKS**

- Existing LLMs can be used to construct questions from a small number of papers (review papers, white papers).
  - Deeper than current MCQs, yet not hyper-specific to papers.
  - Scaling to O(10000) questions is relatively more straightforward.

"question": "If a new particle was discovered that interacts with the weak force but not the electromagnetic or strong force, what could be inferred about its properties?", "distractors": ["It must have a fractional electric charge", "It must be massless", "It could be a candidate for dark matter", "It must be a type of quark"],

"correct\_answer": "It could be a candidate
 for dark matter",



"question": "What is the predominant model
for cosmic ray (CR) transport in the
Galaxy at energies below about 10^17 eV?",
"distractors": [
"Convective wind model with uniform
source distribution",
"Static halo model with isotropic source
emission".
"Leaky-box model with constant escape
time"
1.
"correct answer": "Diffusion model with
possible inclusion of convection".
"question": "If gauge bosons were discovered
"question": "If gauge bosons were discovered to have non-zero rest mass without the
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?",
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the Standard Model to include mass terms for
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the Standard Model to include mass terms for gauge bosons", "The invalidation of gauge
"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the Standard Model to include mass terms for gauge bosons", "The invalidation of gauge symmetry", "The discovery of new gauge
<pre>"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the Standard Model to include mass terms for gauge bosons", "The invalidation of gauge symmetry", "The discovery of new gauge bosons beyond the Standard Model"],</pre>
<pre>"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the Standard Model to include mass terms for gauge bosons", "The invalidation of gauge symmetry", "The discovery of new gauge bosons beyond the Standard Model"], "correct_answer": "A revision of the Standard</pre>
<pre>"question": "If gauge bosons were discovered to have non-zero rest mass without the need for a Higgs mechanism, what would be the most likely implication?", "distractors": ["The existence of a fifth fundamental force", "A revision of the Standard Model to include mass terms for gauge bosons", "The invalidation of gauge symmetry", "The discovery of new gauge bosons beyond the Standard Model"], "correct_answer": "A revision of the Standard Model to include mass terms for gauge</pre>

## **BACKUP: FOUNDATION MODELS BEYOND TEXT: FINE-TUNING**

- Missing features in current LLMs:
  - Support for non-text data like scalar numbers, vectors, images, trees, unstructured datasets.
- Can pre-trained LLMs be re-purposed for numerical tasks?
  - It depends!
  - Degradation/forgetting is expected





- Q: What type of galaxy is expected at redshift of z=0.211? List the telescope surveys where such a galaxy can be observed.
- A galaxy at redshift z=0.211 is expected to be a luminous infrared galaxy (LIRG) or an ultraluminous infrared galaxy (ULIRG) due to the fact that the universe was in its peak star formation period during that epoch. The galaxy is likely to be a merger or a strongly interacting system, and it could be a dusty galaxy with a high rate of star formation.

The galaxy can be observed in the following telescope surveys:

- Sloan Digital Sky Survey (SDSS): The galaxy can be observed in the SDSS-IV survey, which covers the sky in five bands (u, g, r, i, and z) and has a limiting magnitude of 22.5 in the r-band.
- Dark Energy Survey (DES): The galaxy can be observed in the DES survey, which covers the sky in five bands (g, r, i, z, and Y) and has a limiting magnitude of 24.5 in the i-band.
- 3. Hyper Suprime-Cam (HSC) Survey: The galaxy can be observed in the HSC survey, which covers the sky in five bands (g, r, i, z, and Y) and has a limiting magnitude of 26.5 in the i-band.
- 4. Subaru Hyper Suprime-Cam (HSC) Survey: The galaxy can be observed in the Subaru HSC survey, which covers the sky in five bands (g, r, i, z, and Y) and has a limiting magnitude of 26.5 in the i-band.

## **BACKUP: AURORA-GPT — OPEN SCIENCE FOUNDATION MODELS**

- Initially with textual interface that would respond to scientific prompts
- Single model or Mixture/Combination of Experts (scale: ~1 T parameter)





Figure inspired by Bommasani, Rishi, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein et al. "On the opportunities and risks of foundation models." *arXiv preprint arXiv:2108.07258* (2021).