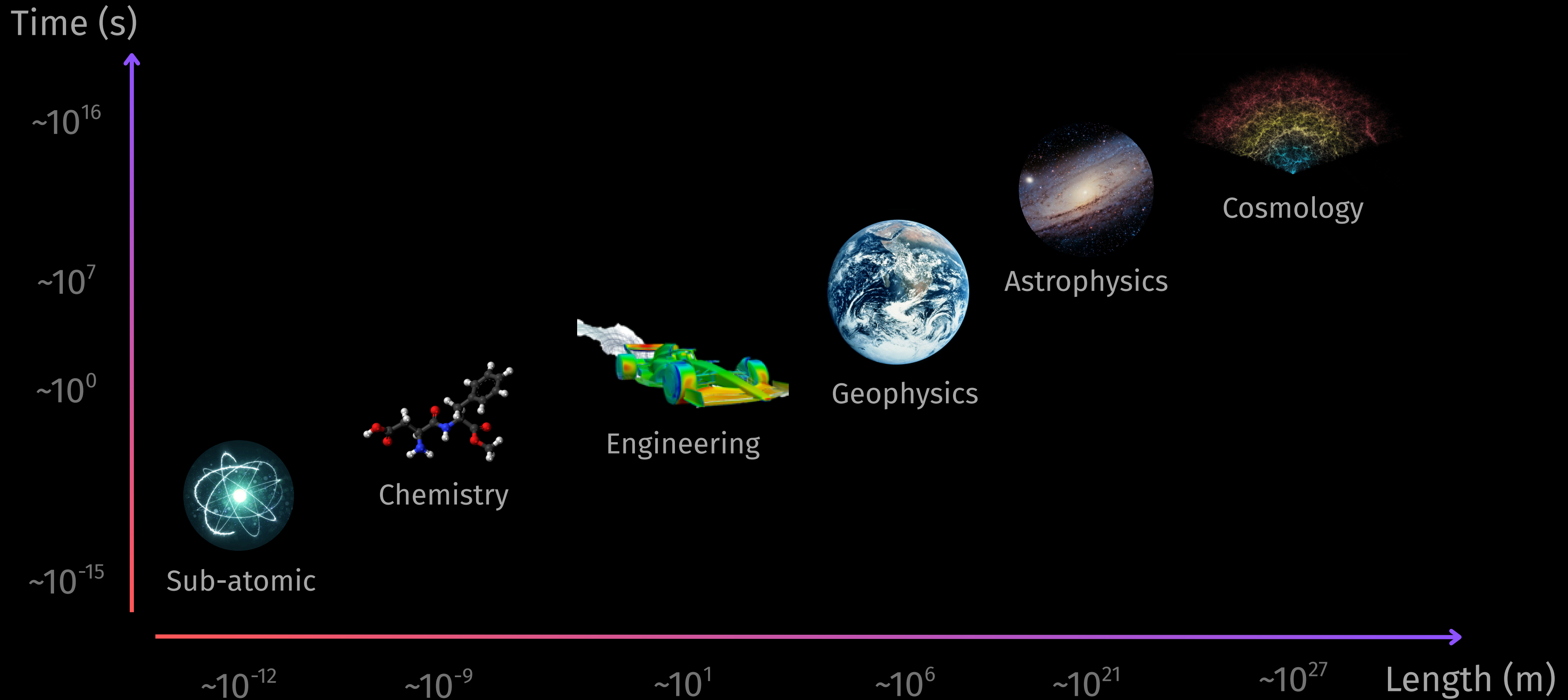

Greener Science with AI: Sustainable Inference of Physical Systems

Jason McEwen
Mission Director for Fundamental Research

Centre for Intelligent Sustainable Computing Symposium
Queen's University Belfast, August 2025



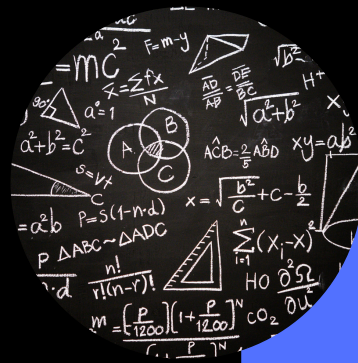
Spatial and temporal scales of physical systems



Pillars of science



1st Pillar:
Experimental



2nd Pillar:
Theoretical



3rd Pillar:
Simulation



4th Pillar:
Data-Driven

~1500

~1700

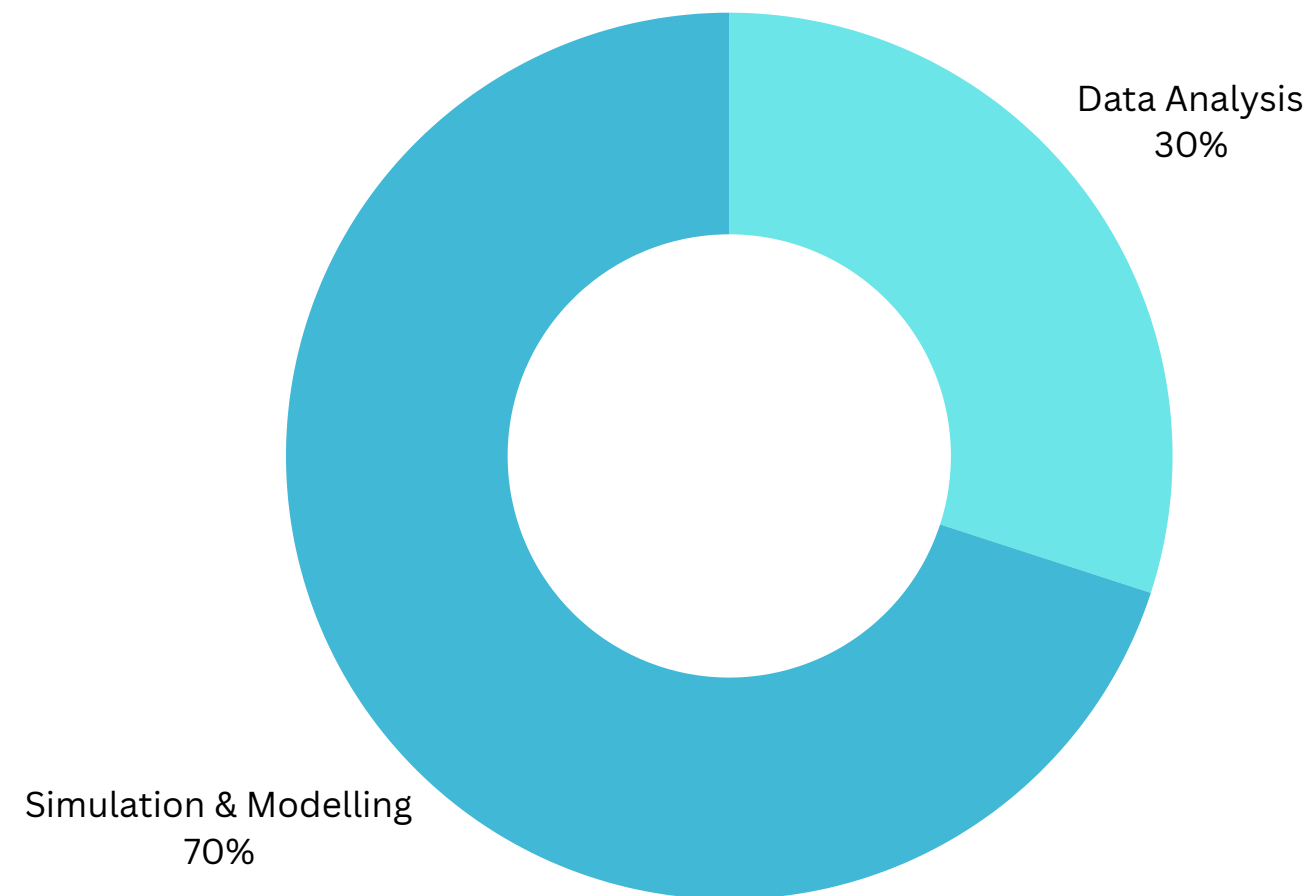
~1950

~2010

t

Computational cost of simulation

- Modelling & simulation account for ~70% of high-performance computing (HPC) usage
- Data analysis accounts for ~30%



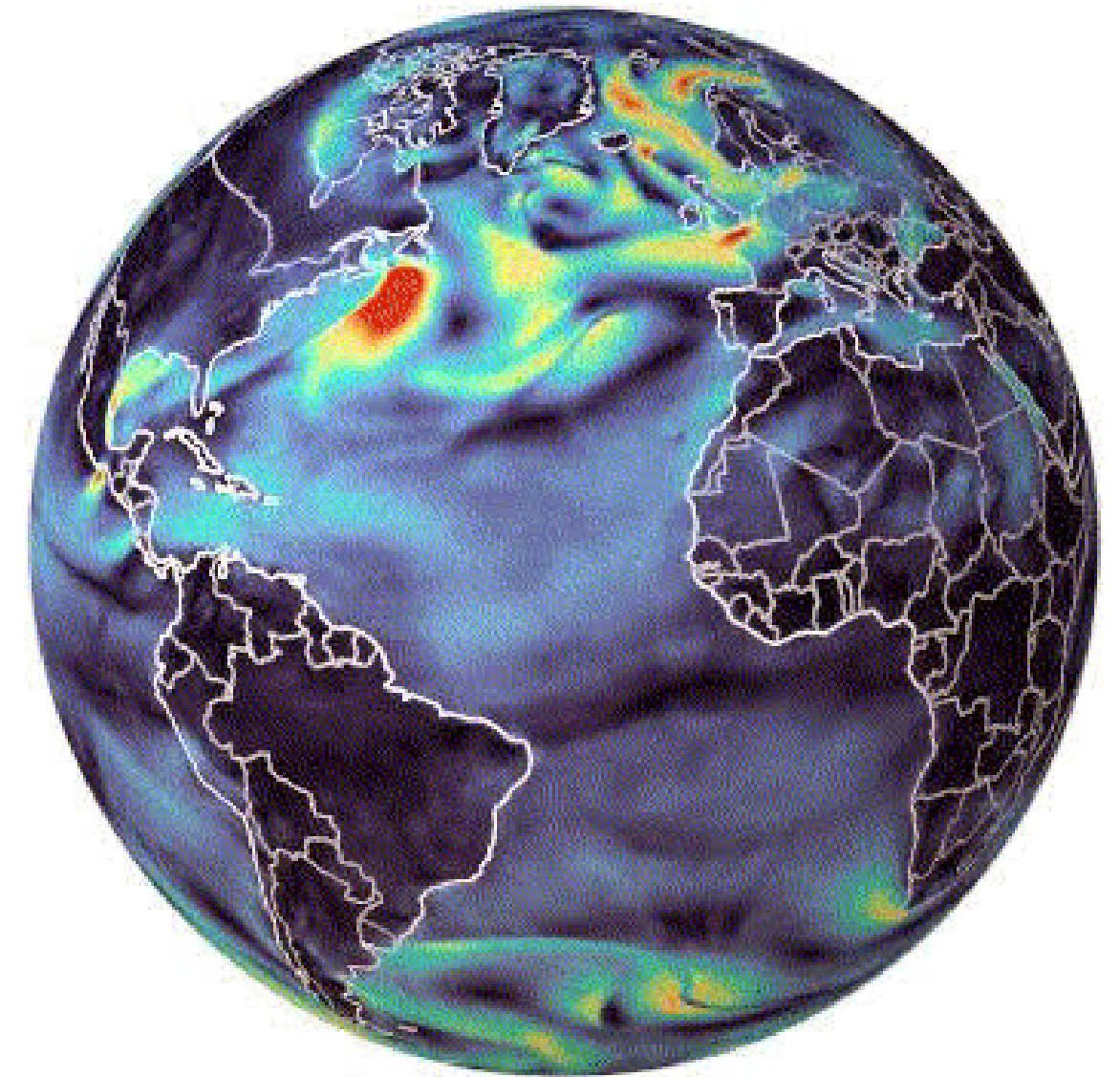
Computational cost of simulation: climate simulation example

CMIP6: Coupled Model Intercomparison Project 6

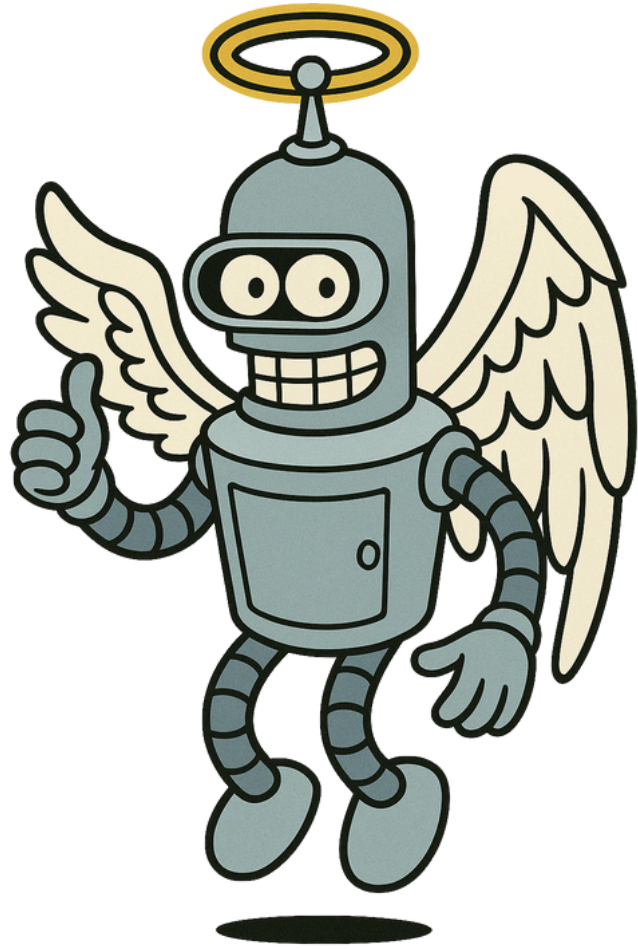
- 40,000 years of climate simulations
- ≥ 1 billion core hours
- 40PB of data
- Carbon footprint of 1692t CO₂ equivalent (Acosta et al. 2024)
 ≈ 6.2 million car miles \approx driving around Earth 250x

Still not enough!

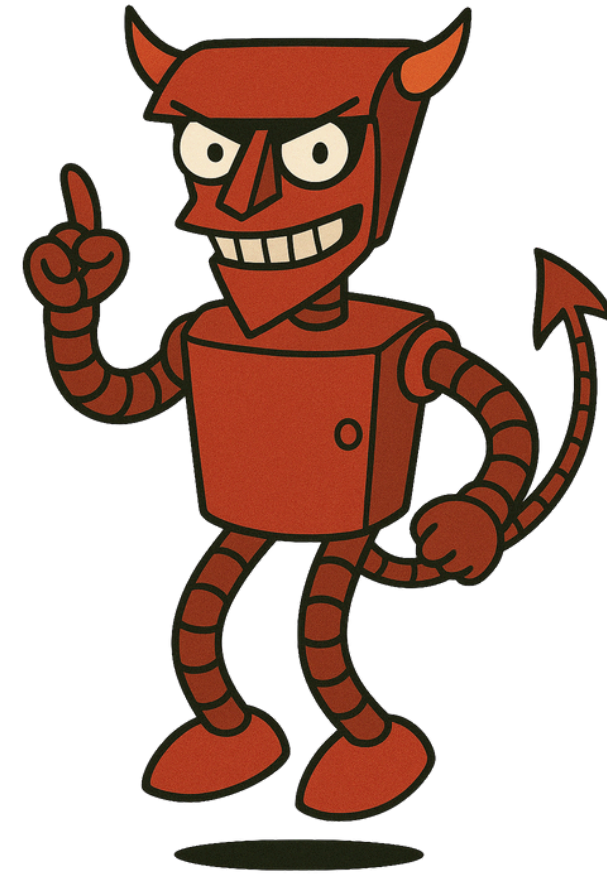
- Need higher resolution
 (currently $\sim 1^\circ \approx 100$ km but require $0.01^\circ \approx 1$ km; Palmer 2014)
- Too few ensembles for robust scenario analysis (uncertainties)
- Too few models (combinatorial explosion of parameters/forcings)



AI to the rescue?



Can we use AI to alleviate
the computation cost of
simulating physical systems?

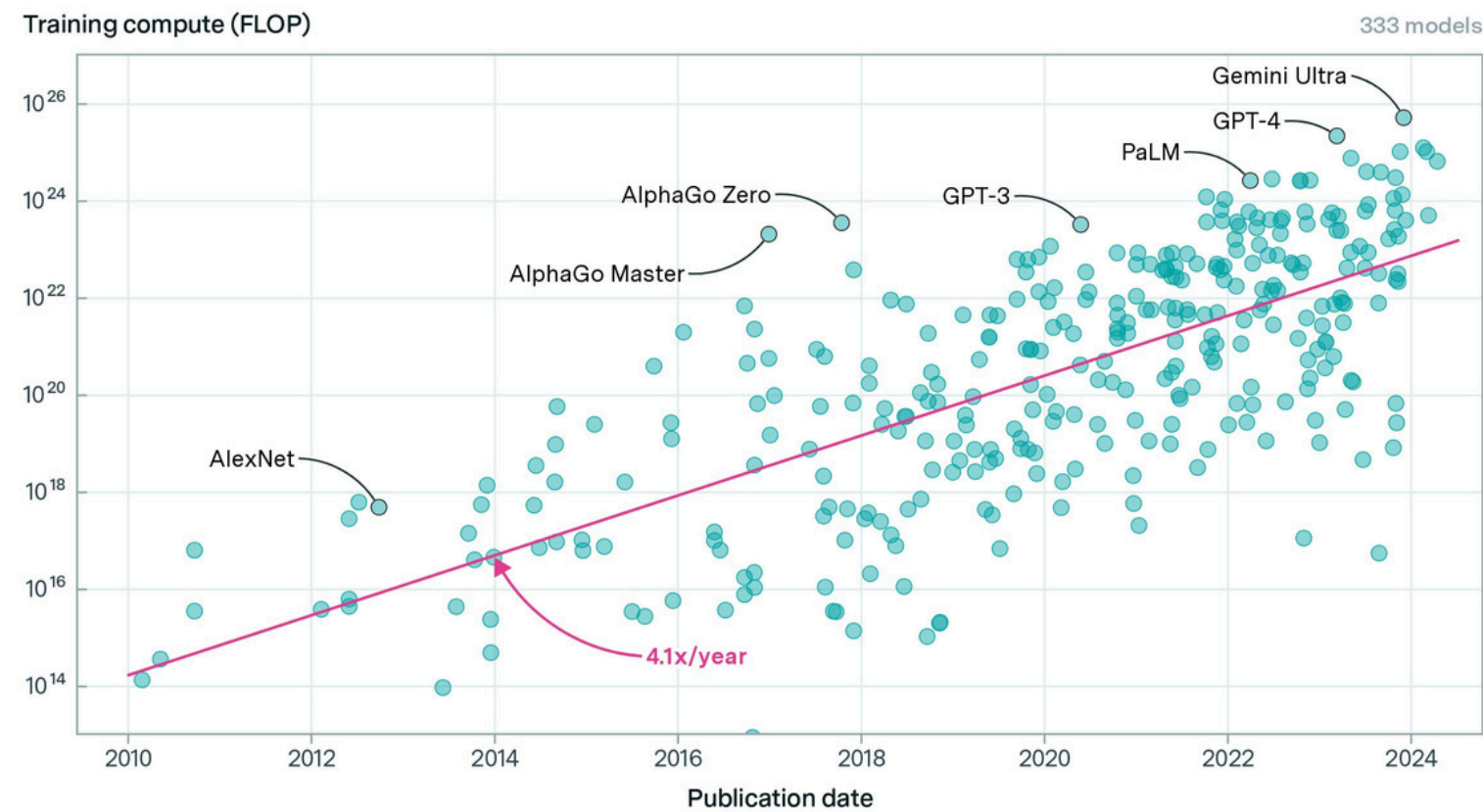


But doesn't AI itself require
huge computational costs?

Computational costs of training large AI models

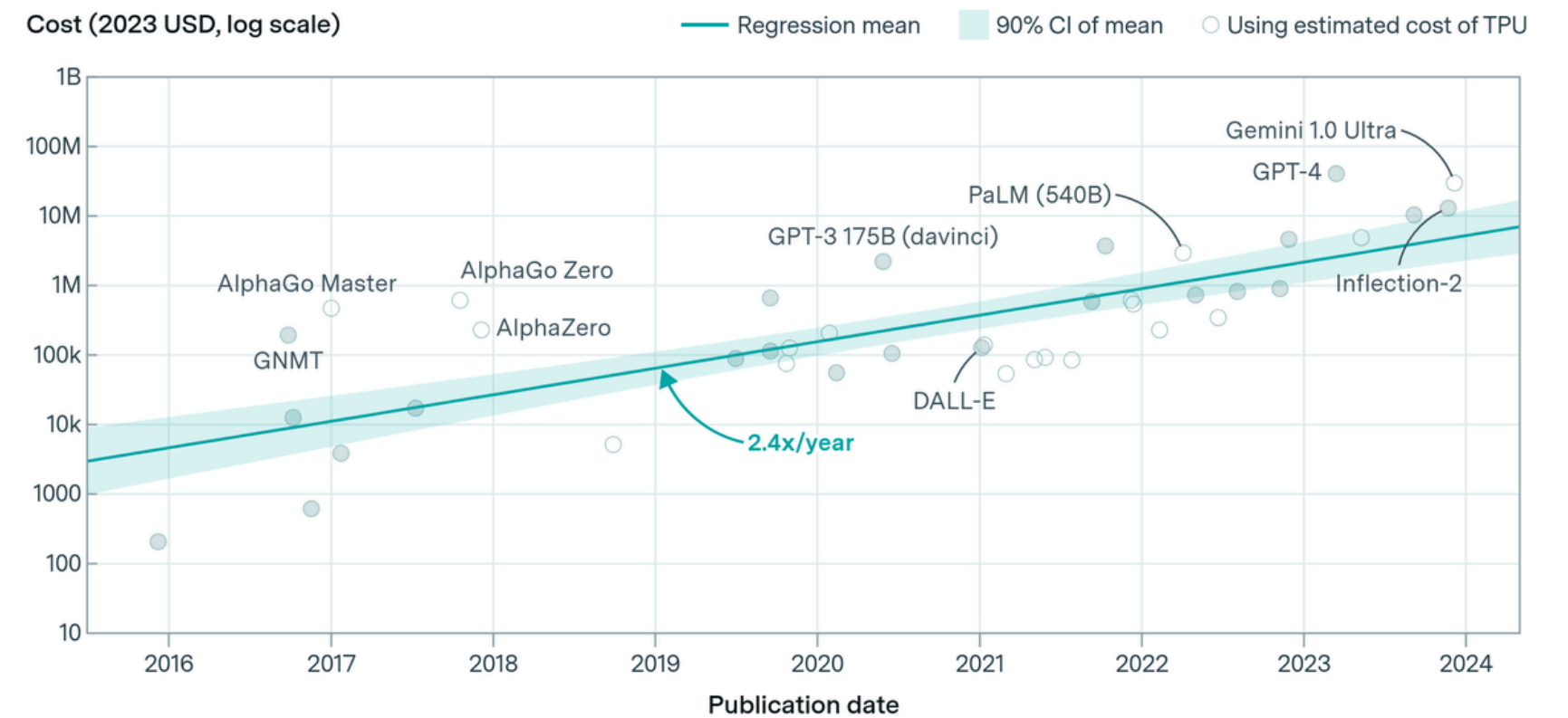
Training compute of notable models

EPOCH AI



Amortized hardware and energy cost to train frontier AI models over time

EPOCH AI



Large vs small AI models



Large AI models

- ✗ Costly (compute, energy, carbon footprint)
- ✓ General purpose

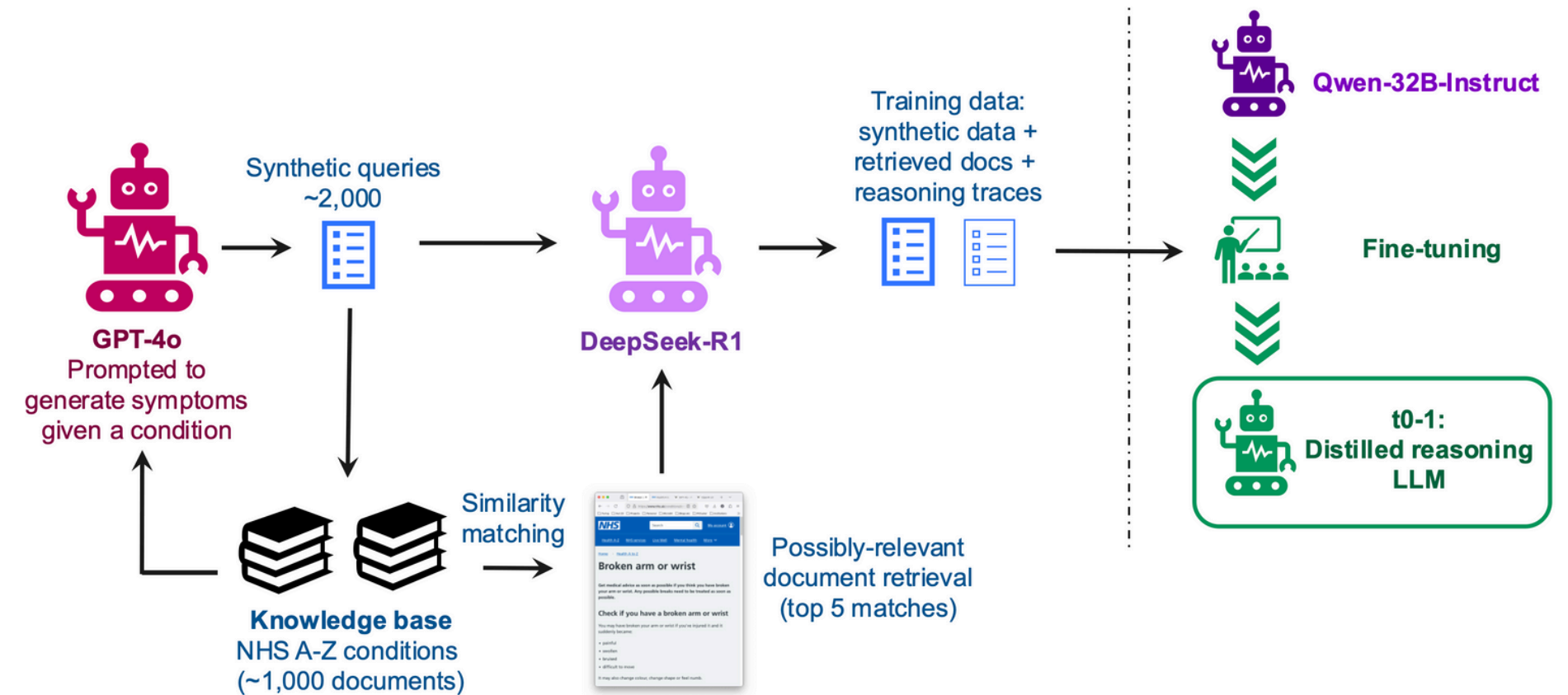


Small AI models

- ✓ Cheap (compute, energy, carbon footprint)
- ✓ Specialised

RAG-augmented reasoning with lean language models (Chan et al. 2025)

- Distillation
 - Fine-tuning
 - Reasoning with budget forcing
 - RAG (retrieval-augmented generation)
- ✓ Frontier performance, without frontier compute
 - ✓ Compute-constrained environments
 - ✓ Privacy-sensitive environments



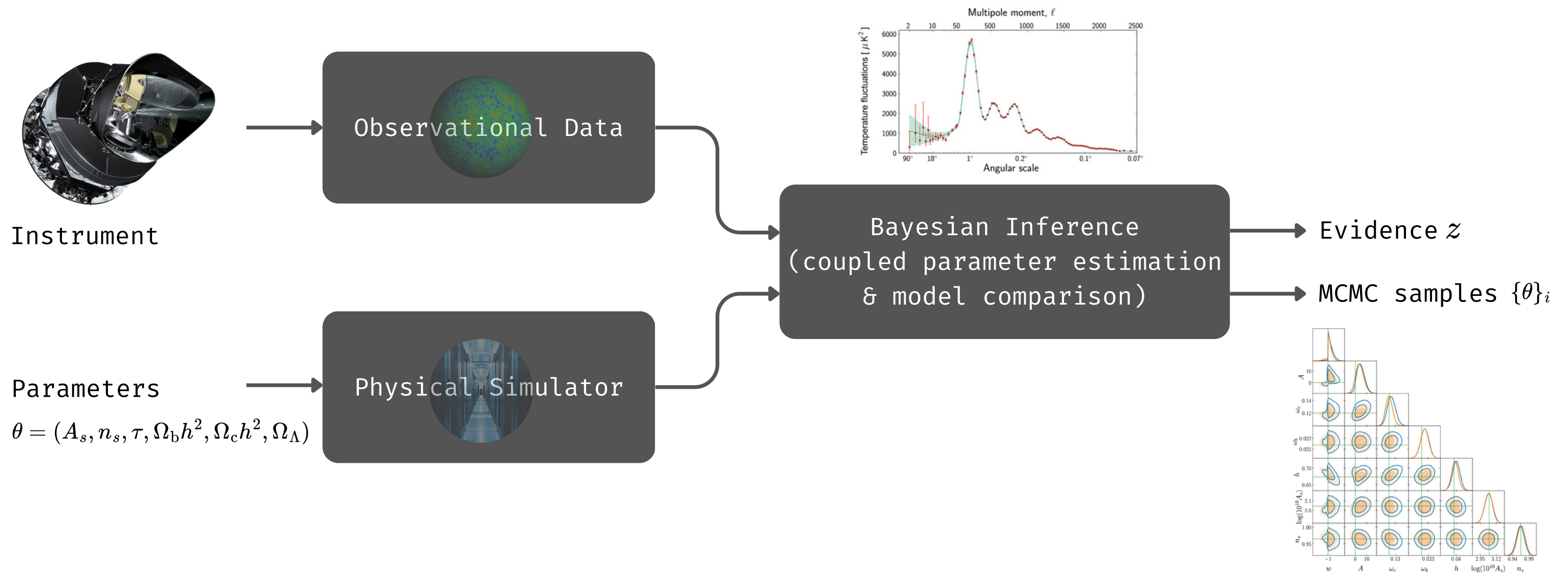
Blog: Why we still need small language models

Outline

1. **Traditional scientific inference** for physical systems
2. **Accelerated scientific inference** for physical systems
 - Emulation
 - Programming frameworks
 - Gradient-based MCMC sampling
 - Decoupled Bayesian model comparison
3. Cosmological **case studies**

Traditional scientific inference for physical systems

Traditional Bayesian inference for physical systems



Bayesian inference: parameter estimation

Bayes' theorem

$$\begin{array}{c} \text{posterior} \\ p(\theta | x, M) \end{array} = \frac{\begin{array}{c} \text{likelihood} \\ p(x | \theta, M) \end{array} \begin{array}{c} \text{prior} \\ p(\theta | M) \end{array}}{\begin{array}{c} p(x | M) \\ \text{marginal likelihood} \end{array}} = \frac{\begin{array}{c} \text{likelihood} \\ \mathcal{L}(\theta) \end{array} \begin{array}{c} \text{prior} \\ \pi(\theta) \end{array}}{\begin{array}{c} z \\ \text{marginal likelihood} \end{array}},$$

for parameters θ , model M and observed data x .

For **parameter estimation**, typically draw samples from the posterior by *Markov chain Monte Carlo (MCMC)* sampling.

Bayesian inference: model comparison

By Bayes' theorem for model M_j :

$$p(M_j | x) = \frac{p(x | M_j)p(M_j)}{\sum_j p(x | M_j)p(M_j)} .$$

For **model comparison**, consider posterior model odds:

$$\underbrace{\frac{p(M_1 | x)}{p(M_2 | x)}}_{\text{posterior odds}} = \underbrace{\frac{p(x | M_1)}{p(x | M_2)}}_{\text{Bayes factor}} \times \underbrace{\frac{p(M_1)}{p(M_2)}}_{\text{prior odds}} .$$

Must compute the **marginal likelihood** (aka. **Bayesian model evidence**) given by the normalising constant

$$z = p(x | M) = \int d\theta \mathcal{L}(\theta) \pi(\theta) .$$

⇒ **Challenging computational problem.**

Nested sampling (Skilling 2006)

Group the parameter space Ω into a series of **nested subspaces**:
 $\Omega_{L^*} = \{x \mid \mathcal{L}(x) \geq L^*\}$. Define the prior volume ξ within Ω_{L^*} by

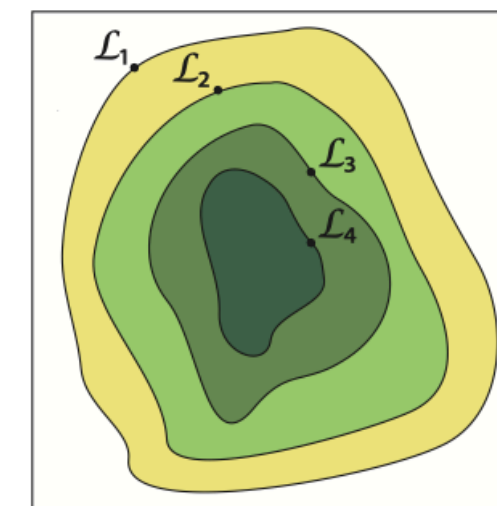
$$\xi(L^*) = \int_{\Omega_{L^*}} \pi(x) dx.$$

Marginal likelihood can then be rewritten as

$$z = \int_0^1 \mathcal{L}(\xi) d\xi.$$

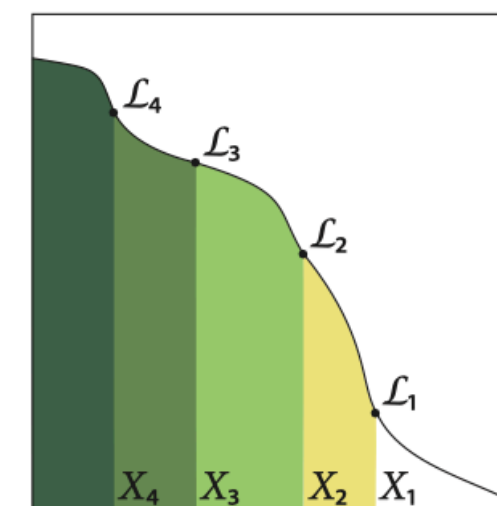
Require computational strategy to compute likelihood level-sets (iso-contours) L_i and corresponding prior volumes $0 < \xi_i \leq 1$.

Crux: sample from the prior, subject to the likelihood level-set constraint, i.e. sample from the prior $\pi(x)$, such that $\mathcal{L}(x) > L^*$.



Feroz et al. (2013)

Nested subspaces

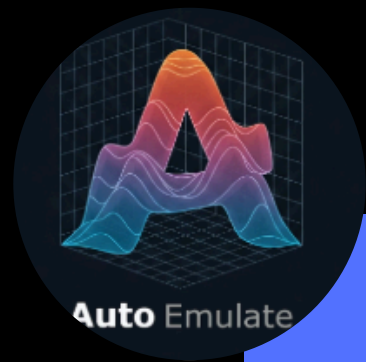


Feroz et al. (2013)

Reparameterised
likelihood

Accelerated scientific inference for physical systems

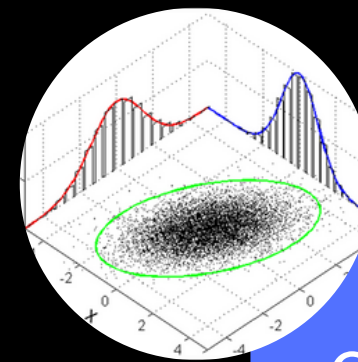
Pillars of accelerated scientific inference



1st Pillar:
AI
Emulation



2nd Pillar:
Programming
Frameworks

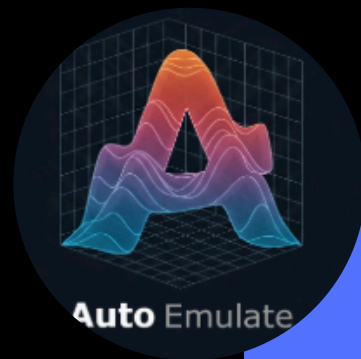


3rd Pillar:
Gradient-Based
MCMC Sampling



4th Pillar:
Decoupled Model
Comparison

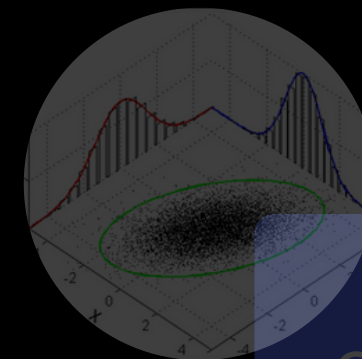
Pillars of accelerated scientific inference



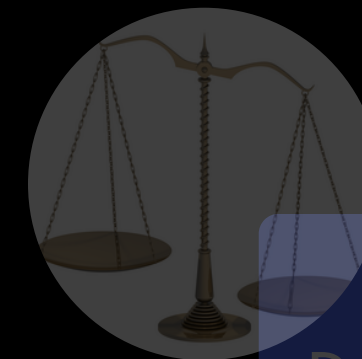
1st Pillar:
AI
Emulation



2nd Pillar:
Programming
Frameworks

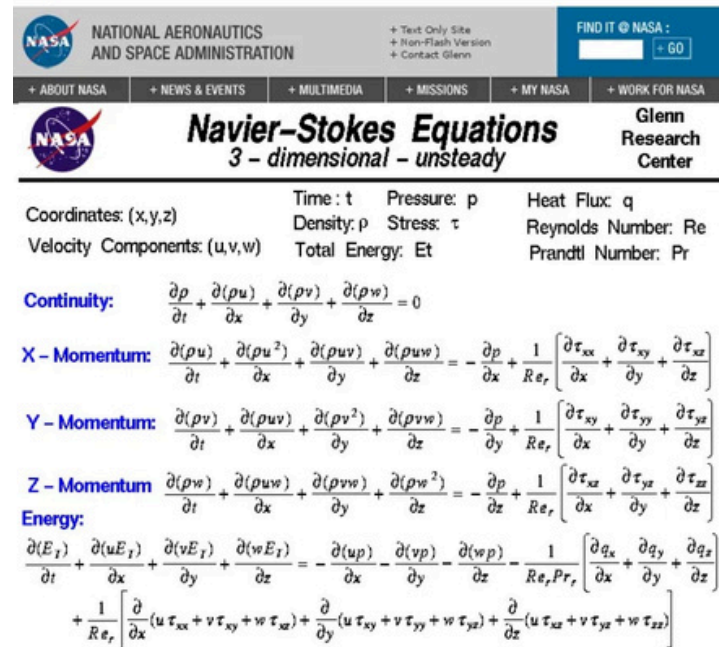


3rd Pillar:
Gradient-Based
MCMC Sampling



4th Pillar:
Decoupled Model
Comparison

Simulation vs emulation



NASA NATIONAL AERONAUTICS AND SPACE ADMINISTRATION

Glenn Research Center

Navier-Stokes Equations
3 - dimensional - unsteady

Coordinates: (x,y,z) Time: t Pressure: p Heat Flux: q
Velocity Components: (u,v,w) Density: ρ Stress: τ Reynolds Number: Re
Total Energy: Et Prandtl Number: Pr

Continuity:
$$\frac{\partial \rho}{\partial t} + \frac{\partial(\rho u)}{\partial x} + \frac{\partial(\rho v)}{\partial y} + \frac{\partial(\rho w)}{\partial z} = 0$$

X - Momentum:
$$\frac{\partial(\rho u)}{\partial t} + \frac{\partial(\rho u^2)}{\partial x} + \frac{\partial(\rho uv)}{\partial y} + \frac{\partial(\rho uw)}{\partial z} = -\frac{\partial p}{\partial x} + \frac{1}{Re_r} \left[\frac{\partial \tau_{xx}}{\partial x} + \frac{\partial \tau_{xy}}{\partial y} + \frac{\partial \tau_{xz}}{\partial z} \right]$$

Y - Momentum:
$$\frac{\partial(\rho v)}{\partial t} + \frac{\partial(\rho uv)}{\partial x} + \frac{\partial(\rho v^2)}{\partial y} + \frac{\partial(\rho vw)}{\partial z} = -\frac{\partial p}{\partial y} + \frac{1}{Re_r} \left[\frac{\partial \tau_{xy}}{\partial x} + \frac{\partial \tau_{yy}}{\partial y} + \frac{\partial \tau_{yz}}{\partial z} \right]$$

Z - Momentum:
$$\frac{\partial(\rho w)}{\partial t} + \frac{\partial(\rho uw)}{\partial x} + \frac{\partial(\rho vw)}{\partial y} + \frac{\partial(\rho w^2)}{\partial z} = -\frac{\partial p}{\partial z} + \frac{1}{Re_r} \left[\frac{\partial \tau_{xz}}{\partial x} + \frac{\partial \tau_{yz}}{\partial y} + \frac{\partial \tau_{zz}}{\partial z} \right]$$

Energy:
$$\frac{\partial(E_t)}{\partial t} + \frac{\partial(uE_t)}{\partial x} + \frac{\partial(vE_t)}{\partial y} + \frac{\partial(wE_t)}{\partial z} = -\frac{\partial(u p)}{\partial x} - \frac{\partial(v p)}{\partial y} - \frac{\partial(w p)}{\partial z} - \frac{1}{Re_r Pr_r} \left[\frac{\partial q_x}{\partial x} + \frac{\partial q_y}{\partial y} + \frac{\partial q_z}{\partial z} \right] + \frac{1}{Re_r} \left[\frac{\partial}{\partial x} (u \tau_{xx} + v \tau_{xy} + w \tau_{xz}) + \frac{\partial}{\partial y} (u \tau_{xy} + v \tau_{yy} + w \tau_{yz}) + \frac{\partial}{\partial z} (u \tau_{xz} + v \tau_{yz} + w \tau_{zz}) \right]$$

Simulate physical laws

- ✓ Accurate representation of physical model
- ✗ Highly computationally costly



Emulate by training an AI model
to micmic physical laws

- ✗ Approximate representation of physical model
- ✓ Computationally efficient (once trained)
- ✓ Learning data-driven model has potential to be more accurate than physical model



Users with Domain expertise



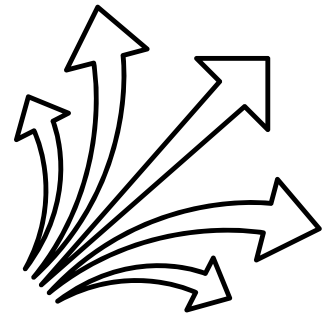
No Machine learning
expertise in Emulation is
necessary



Democratising the use of AI
for accelerating simulations
in various industries



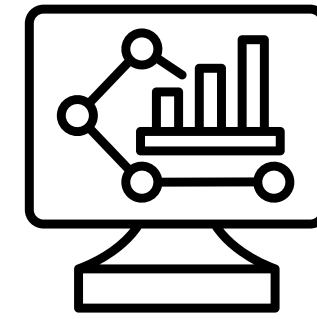
AutoEmulate: more than a learned emulator



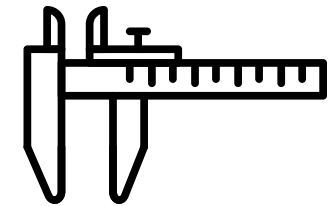
Sensitivity analysis



History matching



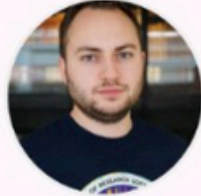














Simulator in the loop

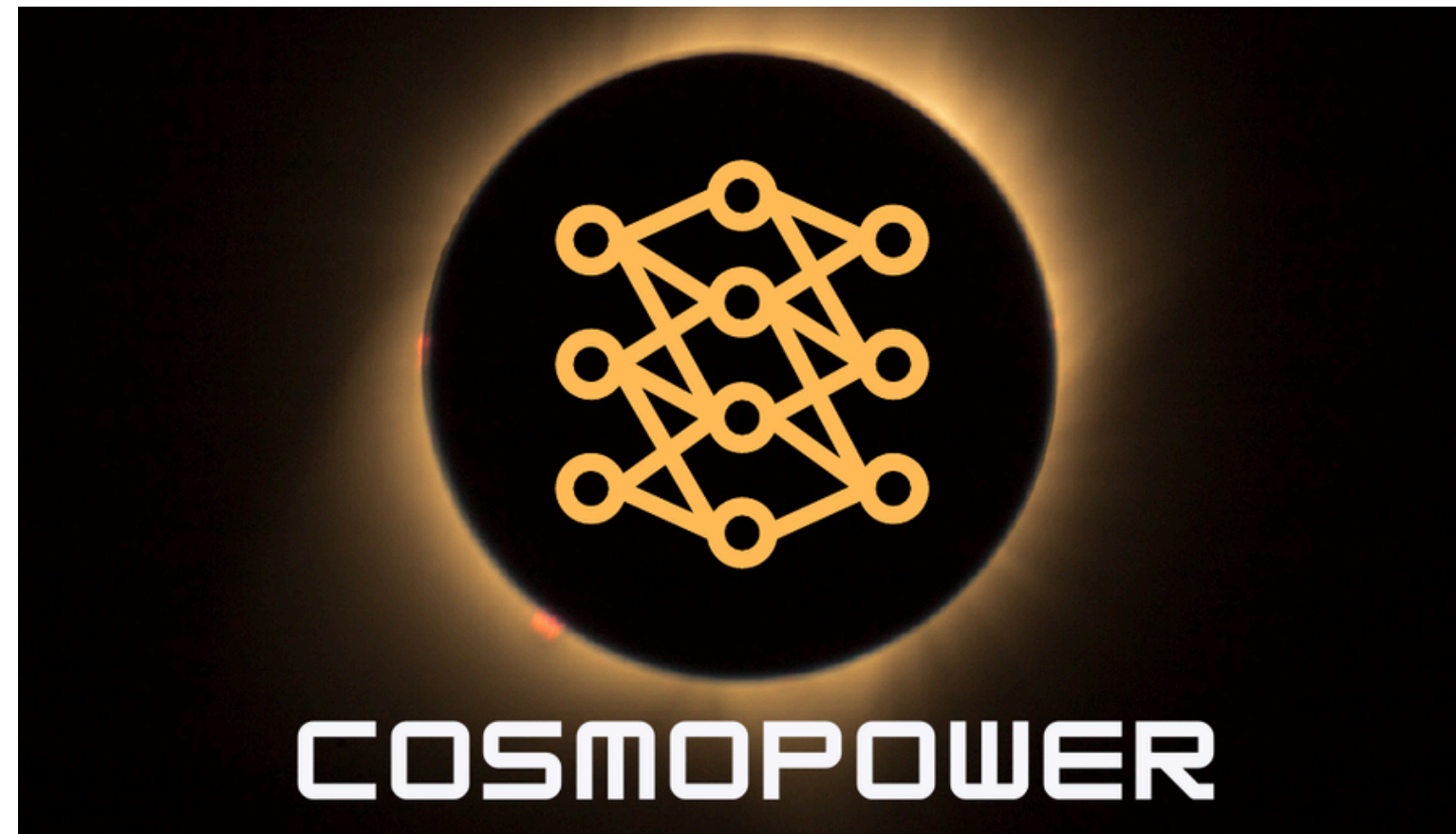


Bayesian calibration

AutoEmulate team

 <p>Andrew Duncan Associate Professor in Statistical Machine Learning at Imperial College London</p>	 <p>Christopher Iliffe Sprague Research Associate</p>	 <p>Ed Chalstrey Research Engineer</p>	 <p>Nayara Fonseca Research Associate</p>	 <p>Paolo Conti Research Associate</p>
 <p>Edwin Brown Research Engineer</p>	 <p>Jason McEwen Mission Director - Fundamental Research in AI for Physical Systems</p>	 <p>Karen de Cesare Research Project Manager</p>	 <p>Sam Greenbury Research Engineer</p>	 <p>Steven Niederer Chair in Biomedical Engineering at Imperial College London</p>
 <p>Marjan Famili Research Associate</p>	 <p>Martin Stoffel Research Engineer</p>	 <p>Maya Bronfeld Programme Manager</p>	 <p>Radka Jersakova Research Engineer</p>	 <p>Tomas Lazauskas Research Computing Team Lead</p>

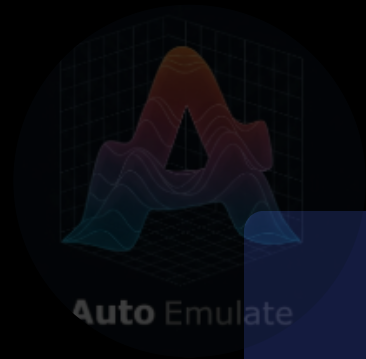
Emulation for cosmology



<https://github.com/alessiospuriomancini/cosmopower>

<https://github.com/dpiras/cosmopower-jax>

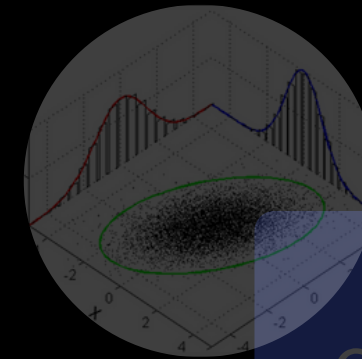
Pillars of accelerated scientific inference



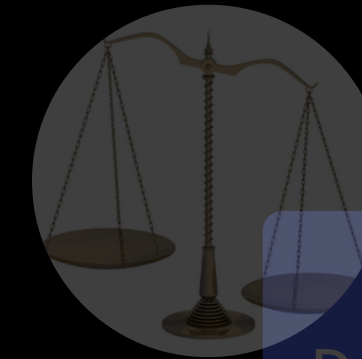
1st Pillar:
AI
Emulation



2nd Pillar:
Programming
Frameworks

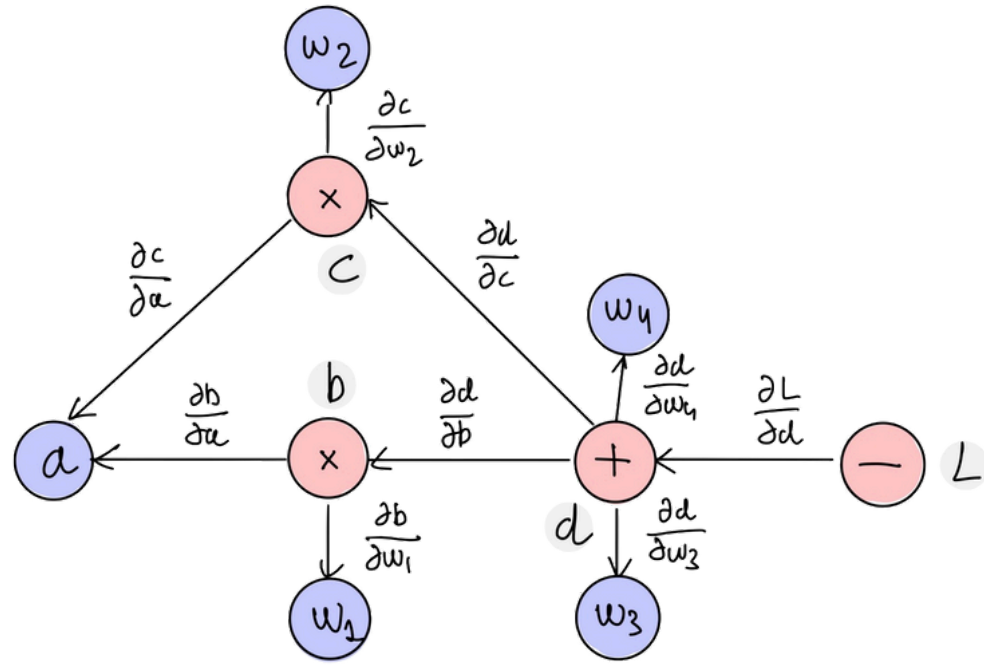


3rd Pillar:
Gradient-Based
MCMC Sampling

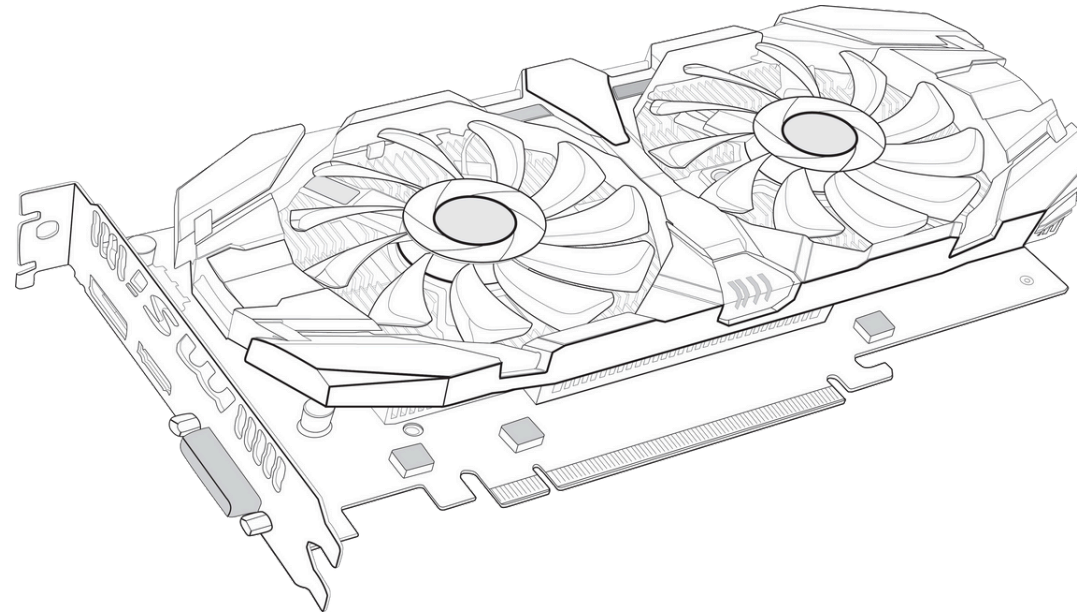


4th Pillar:
Decoupled Model
Comparison

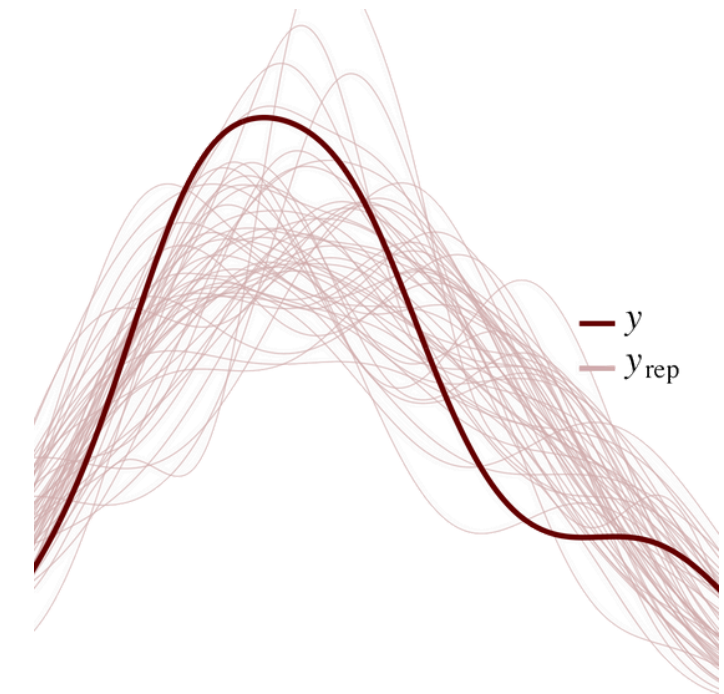
Programming frameworks



Automatic differentiation



GPU acceleration



Probabilistic programming



JAX



XLA



NumPyro

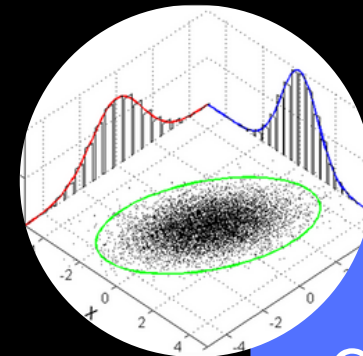
Pillars of accelerated scientific inference



1st Pillar:
AI
Emulation



2nd Pillar:
Programming
Frameworks



3rd Pillar:
Gradient-Based
MCMC Sampling



4th Pillar:
Decoupled Model
Comparison

Gradient-accelerated MCMC sampling

Exploit gradient information to scale MCMC efficiently to higher dimensional settings (e.g. Hamiltonian or Langevin dynamics).

Consider **Hamiltonian Monte Carlo (HMC)**, where samples θ augmented with momentum p . Hamiltonian given by

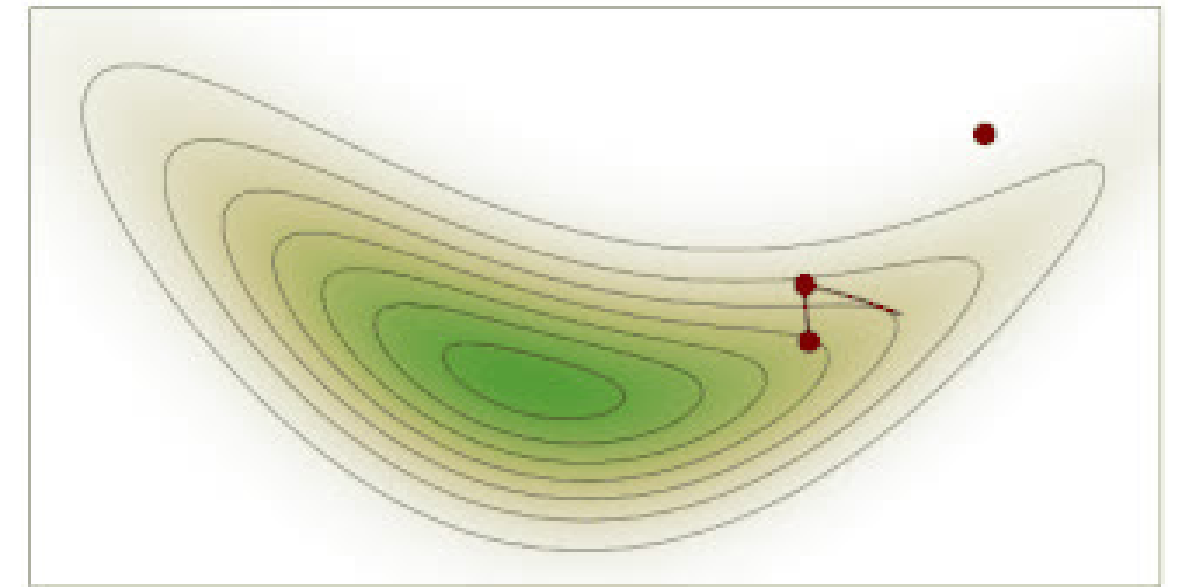
$$H(\theta, p) = -\log p(\theta|x) + \frac{1}{2} p^T M^{-1} p,$$

where M is the mass matrix. Evolution given by dynamics

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial p}, \quad \frac{dp}{dt} = -\frac{\partial H}{\partial \theta}.$$

Consider **No U-Turn (NUTS)** algorithm.

Compute gradients efficiently by automatic differentiation.



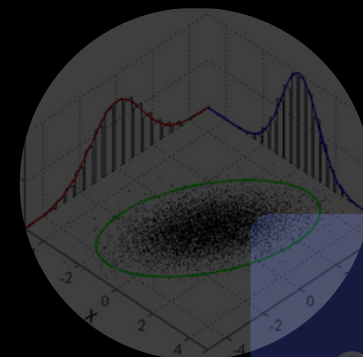
Pillars of accelerated scientific inference



1st Pillar:
AI
Emulation



2nd Pillar:
Programming
Frameworks



3rd Pillar:
Gradient-Based
MCMC Sampling



4th Pillar:
Decoupled Model
Comparison

The problem of nested sampling for Bayesian model comparison

Nested sampling (Skilling 2006) has been the method of choice for almost two decades!

Many highly effective nested sampling algorithms (for a review see Ashton *et al.* 2022).

However, nested sampling has a **fundamental problem**...

Nested sampling tightly couples sampling strategy to marginal likelihood calculation.

As the name suggests, **one must sample in a nested manner.**

- ▷ **Precludes** many alternative **accelerated sampling** strategies that scale to high-dimensions.
- ▷ **Precludes** use in many **simulation-based inference (SBI)** and **variational inference (VI)** settings, where one draws posterior samples directly.

Original harmonic mean estimator

Harmonic mean relationship (Newton & Raftery 1994)

$$\rho = \mathbb{E}_{p(\theta | x)} \left[\frac{1}{\mathcal{L}(\theta)} \right] = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | x) = \int d\theta \frac{1}{\mathcal{L}(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{z} = \frac{1}{z}$$

Original harmonic mean estimator (Newton & Raftery 1994)

$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\mathcal{L}(\theta_i)}, \quad \theta_i \sim p(\theta | x)$$

✓ Only requires posterior samples!

✗ But can fail catastrophically! (Neal 1994)

Importance sampling interpretation of harmonic mean estimator

Alternative interpretation of harmonic mean relationship:

$$\rho = \int d\theta \frac{1}{\mathcal{L}(\theta)} p(\theta | x) = \frac{1}{Z} \overset{\text{importance sampling}}{\int d\theta \frac{\pi(\theta)}{p(\theta | x)} p(\theta | x)} .$$

Importance sampling interpretation:

- ▷ Importance **sampling target distribution** is prior $\pi(\theta)$.
- ▷ Importance **sampling density** is posterior $p(\theta | x)$.

For importance sampling, want sampling density to have fatter tails than target.

Importance sampling failure mode when sampling density is posterior and target is prior.

Re-targeted harmonic mean estimator

Re-targeted harmonic mean relationship (Gelfand & Dey 1994)

$$\rho = \mathbb{E}_{p(\theta | x)} \left[\frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right] = \frac{1}{Z}$$

Normalised distribution $\varphi(\theta)$ now plays the role of the importance sampling target

↪ must **not** have fatter tails than posterior.

Re-targeted harmonic mean estimator (Gelfand & Dey 1994)

$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{\varphi(\theta_i)}{\mathcal{L}(\theta_i)\pi(\theta_i)}, \quad \theta_i \sim p(\theta | x)$$

↪ How set importance sampling target distribution $\varphi(\theta)$?

How set importance sampling target distribution?

Variety of cases been considered:

- ▷ Multi-variate Gaussian (*e.g.* Chib 1995)
- ▷ Indicator functions (*e.g.* Robert & Wraith 2009, van Haasteren 2009)

Optimal target: (McEwen *et al.* 2021)

$$\varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{z}.$$

But clearly **not feasible** since requires knowledge of the evidence z (recall the target must be normalised) \rightsquigarrow **requires problem to have been solved already!**

Learned harmonic mean estimator

Learn an approximation of the optimal target distribution:

$$\varphi(\theta) \stackrel{\text{AI}}{\simeq} \varphi^{\text{optimal}}(\theta) = \frac{\mathcal{L}(\theta)\pi(\theta)}{Z} .$$

- ▷ Approximation not required to be highly accurate.
- ▷ Critically, **must not have fatter tails than posterior.**

Constraining tails of target approach 1: bespoke optimisation problem

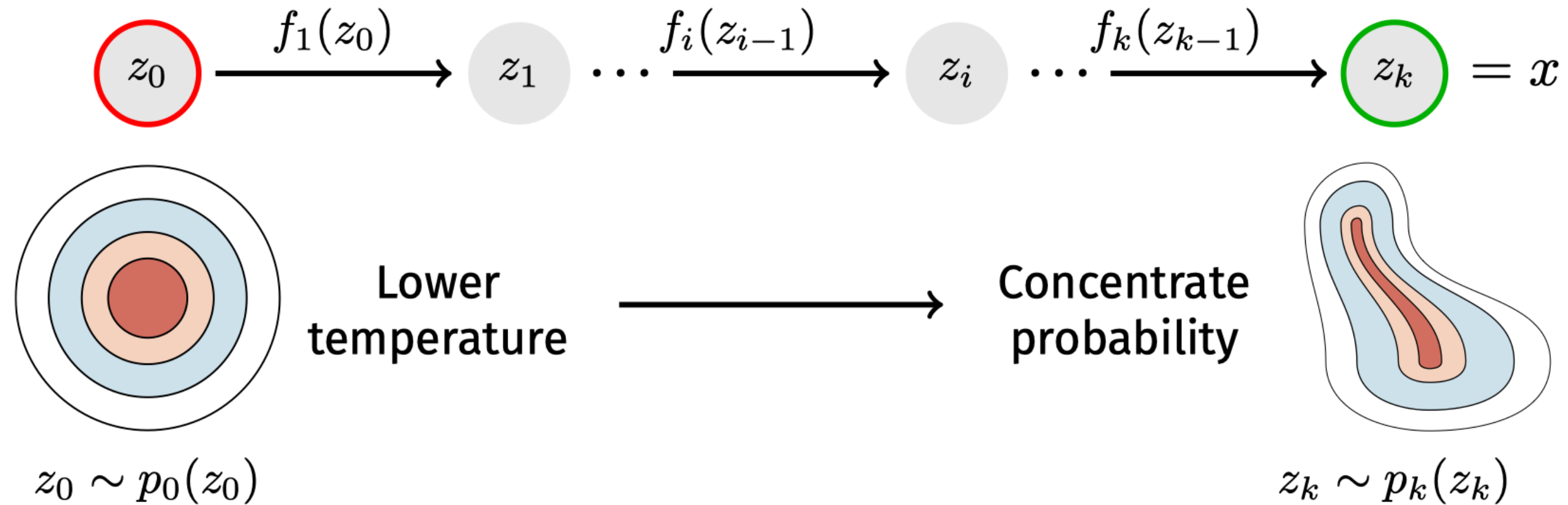
Fit density estimator by **minimising variance of resulting estimator**, with possible regularisation:

$$\min \hat{\sigma}^2 + \lambda R \quad \text{subject to} \quad \hat{\rho} = \hat{\mu}_1.$$

Solve by bespoke **mini-batch stochastic gradient descent**.

Cross-validation to select density estimation model and hyperparameters.

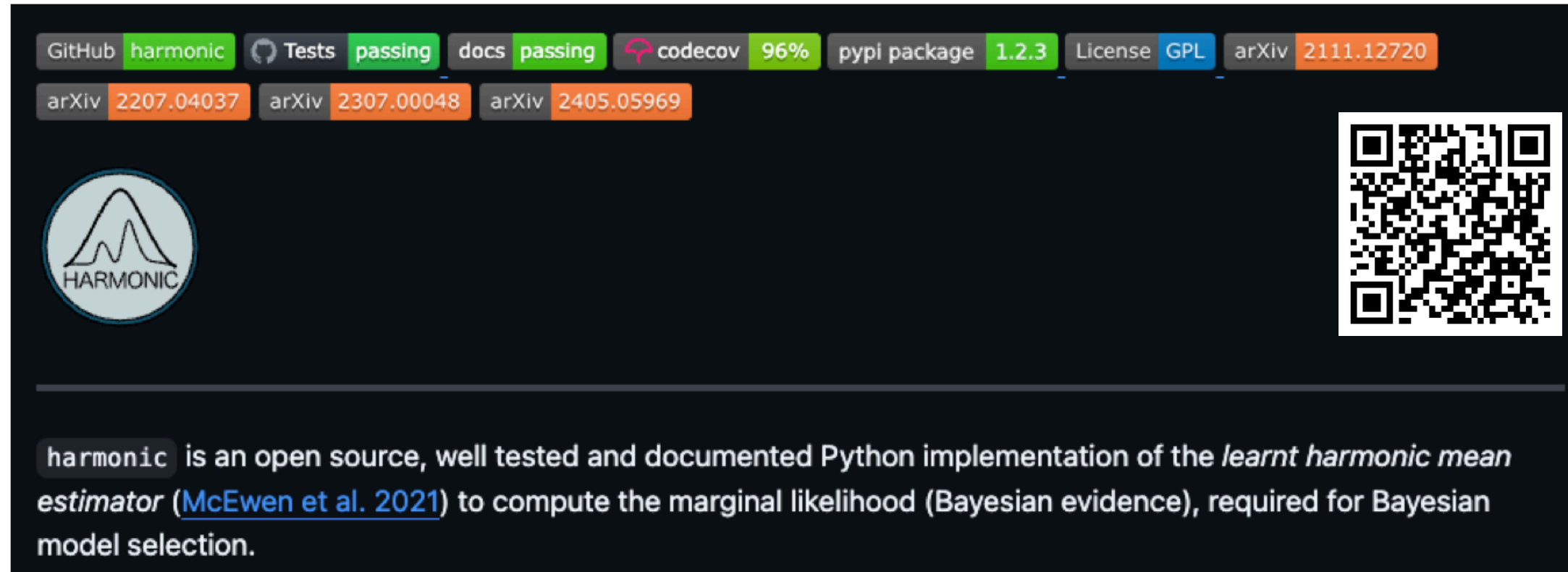
Constraining tails of target approach 2: normalizing flows



- ✓ **Flexible:** no bespoke training; can vary T after training.
- ✓ **Robust:** only one hyperparameter T that does not require fine tuning.
- ✓ **Scalable:** flows scale to higher dimensions than classical density estimators.

(Polanska et al. McEwen 2024)

Harmonic code

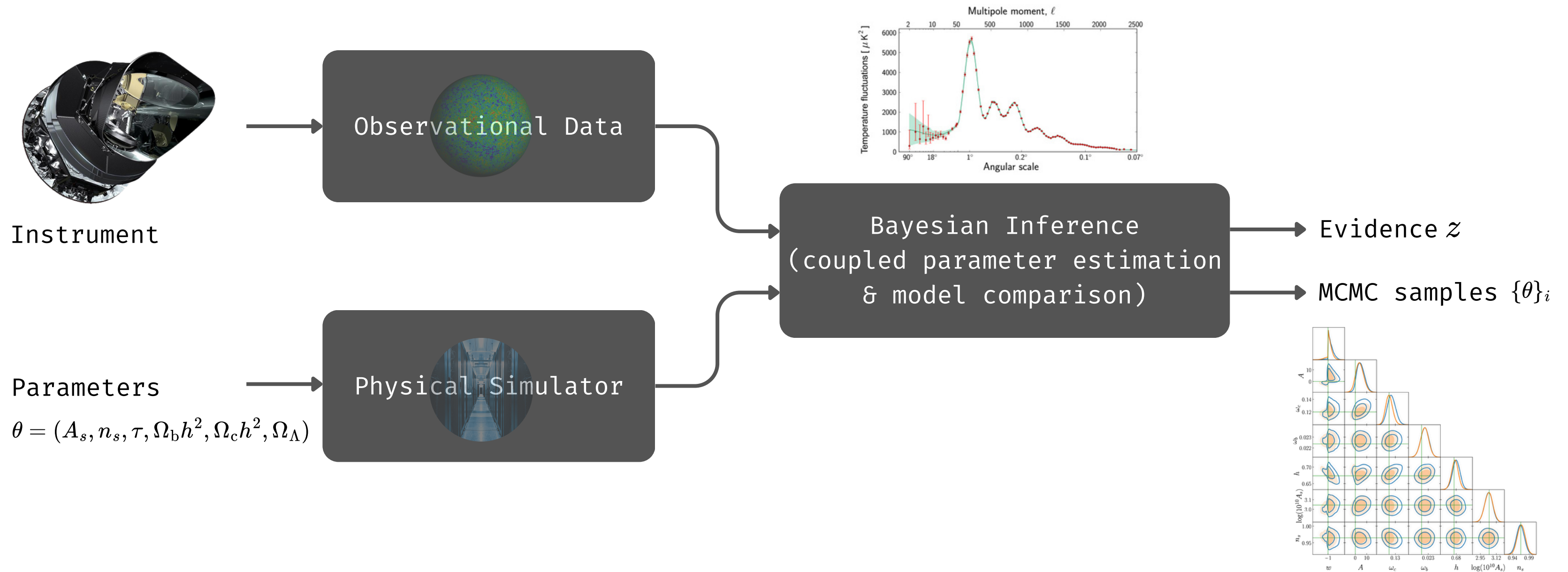


The screenshot displays the project page for 'harmonic'. At the top, there are several badges: GitHub (harmonic), Tests (passing), docs (passing), codecov (96%), pypi package (1.2.3), License (GPL), and arXiv (2111.12720). Below these are three more arXiv badges with IDs 2207.04037, 2307.00048, and 2405.05969. On the left is a circular logo with a mountain-like shape and the word 'HARMONIC' below it. On the right is a QR code. A horizontal line separates the header from the description. The description states: 'harmonic is an open source, well tested and documented Python implementation of the *learnt harmonic mean estimator* (McEwen et al. 2021) to compute the marginal likelihood (Bayesian evidence), required for Bayesian model selection.'

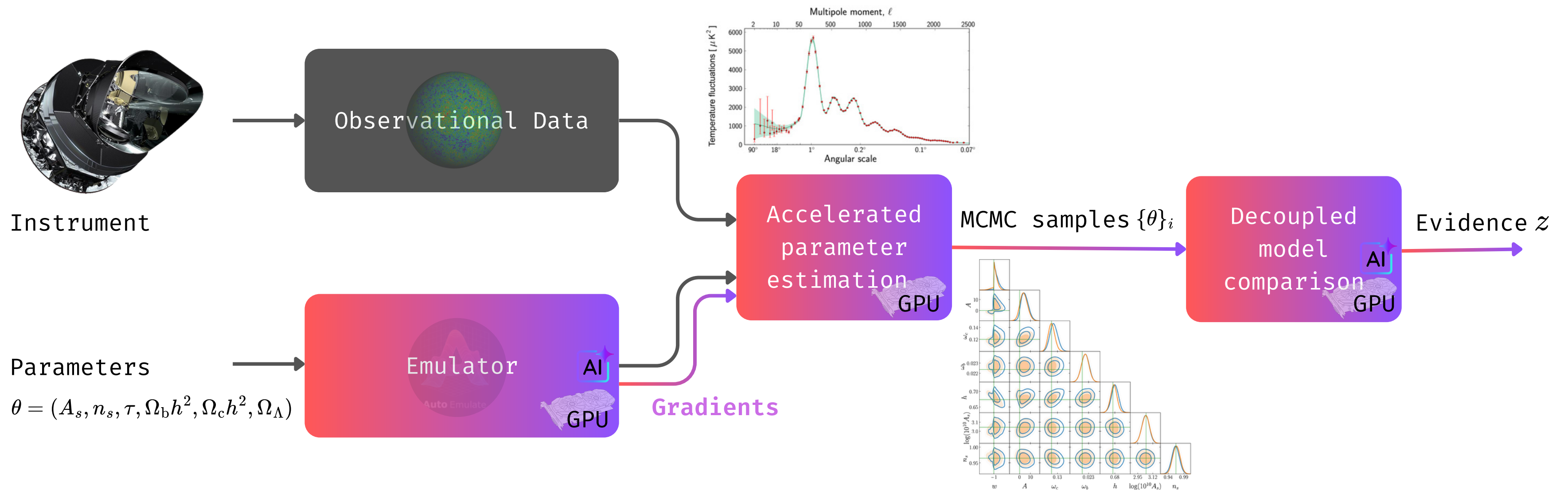
Github: <https://github.com/astro-informatics/harmonic>

Docs: <https://astro-informatics.github.io/harmonic>

Traditional Bayesian inference for physical systems



Accelerated Bayesian inference for physical systems



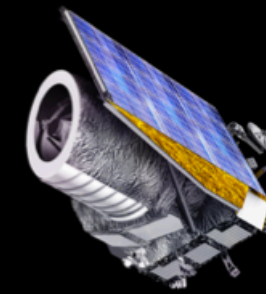
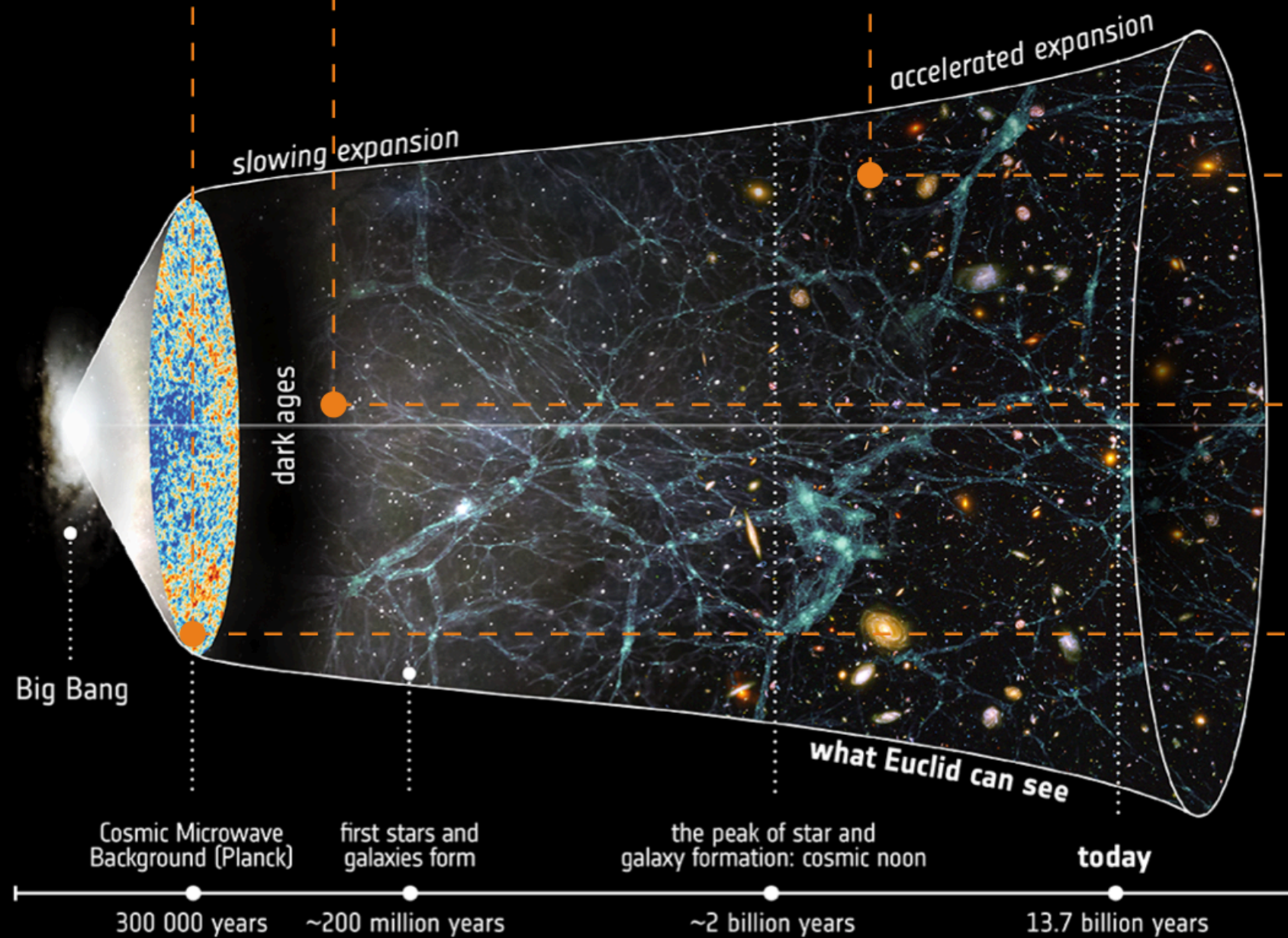
Cosmological **case studies**

Towards a fundamental understanding of our Universe

What is the origin of structure?

How did luminous large-scale structure form?

What is the nature of dark energy and dark matter?



Euclid



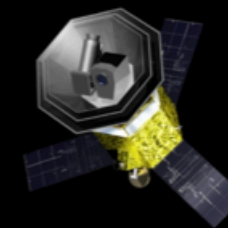
Roman



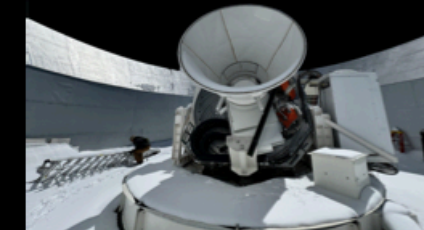
Rubin-LSST



SKA



LiteBIRD



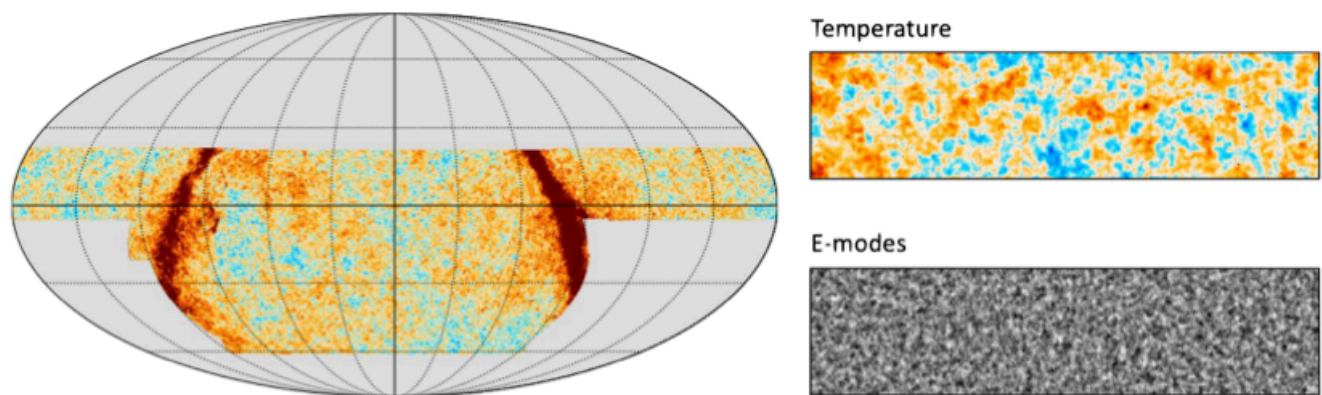
Simons

Atacama Cosmology Telescope (ACT) analysis

Compare Λ CDM (Einstein’s cosmological constant) vs w_0w_a CDM (dynamical dark energy) using learned harmonic mean (McEwen *et al.*2021) with ACT data (Aiola *et al.* 2020).



Atacama Cosmology Telescope (ACT)



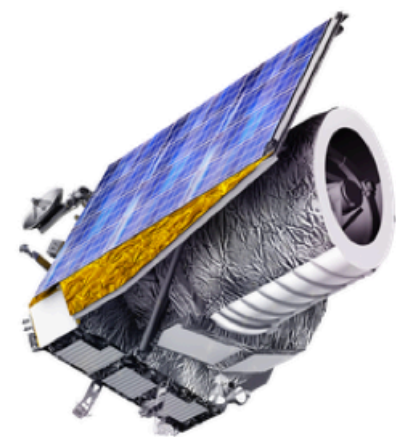
CMB observations

7D vs 9D models:	Λ CDM	w_0w_a CDM	$\log \text{BF}_{\Lambda\text{CDM}-w_0w_a\text{CDM}}$
Nested sampling	-168.92 ± 0.35	-169.38 ± 0.24	0.46 ± 0.42
Learned harmonic mean	-168.87 ± 0.29	-169.32 ± 0.25	0.45 ± 0.38

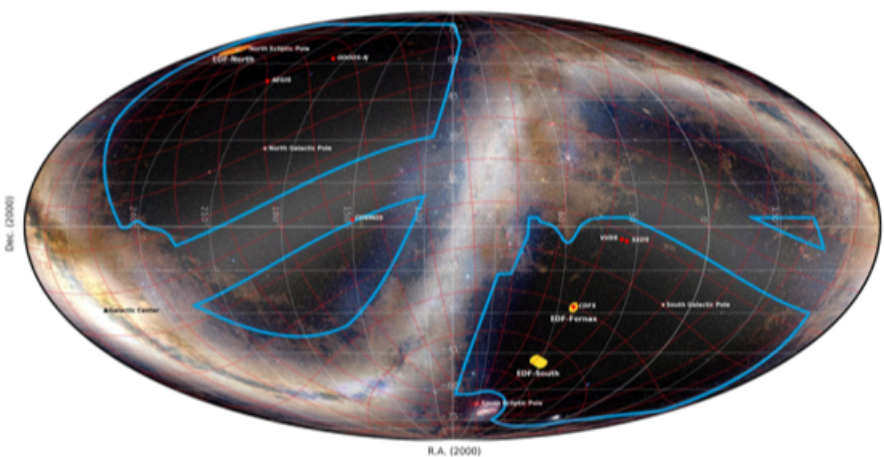
$\rightsquigarrow \Lambda$ CDM mildly favoured \rightsquigarrow **3× acceleration** (Only Pillar 4)

Euclid (Stage IV survey)-like analysis

Compare Λ CDM vs w_0w_a CDM leveraging **4 pillars of AI-acceleration** with Euclid-like lensing and clustering simulations (Piras *et al.* 2024).



Euclid satellite



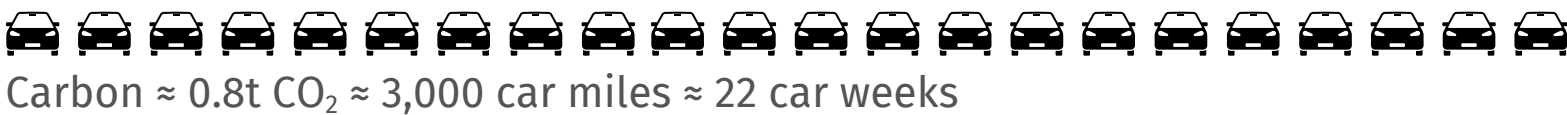
Observation field

37D vs 39D models:	$\log(z_{\Lambda\text{CDM}})$	$\log(z_{w_0w_a\text{CDM}})$	$\log \text{BF}_{\Lambda\text{CDM}-w_0w_a\text{CDM}}$	Total computation time
Classical	-107.03 ± 0.27	-107.81 ± 0.74	0.78 ± 0.79	8 months (48 CPUs)
AI-accelerated (ours)	40956.55 ± 0.06	40955.03 ± 0.04	1.53 ± 0.07	2 days (12 GPUs)

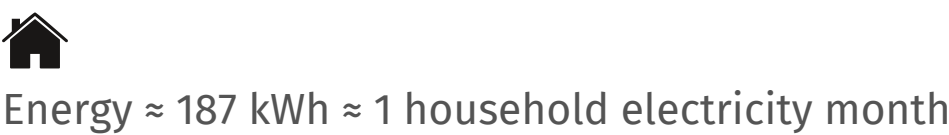
Simulating training data = 1 CPU day | Training = 1 GPU hour | Amortized over all analyses

Euclid (Stage IV survey)-like analysis

Traditional approach

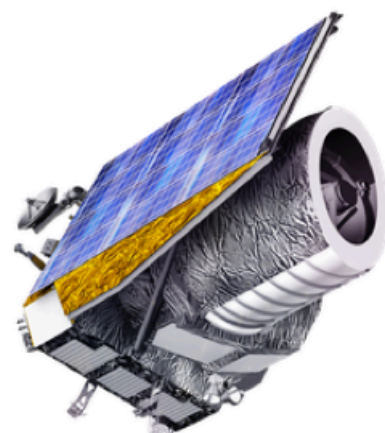


Accelerated approach (ours)



Euclid-Rubin-Roman (3x Stage IV survey)-like analysis

Extend to combined 3× Stage IV Survey-like lensing and clustering simulations (Piras *et al.* 2024).



Euclid satellite



Rubin observatory



Roman satellite

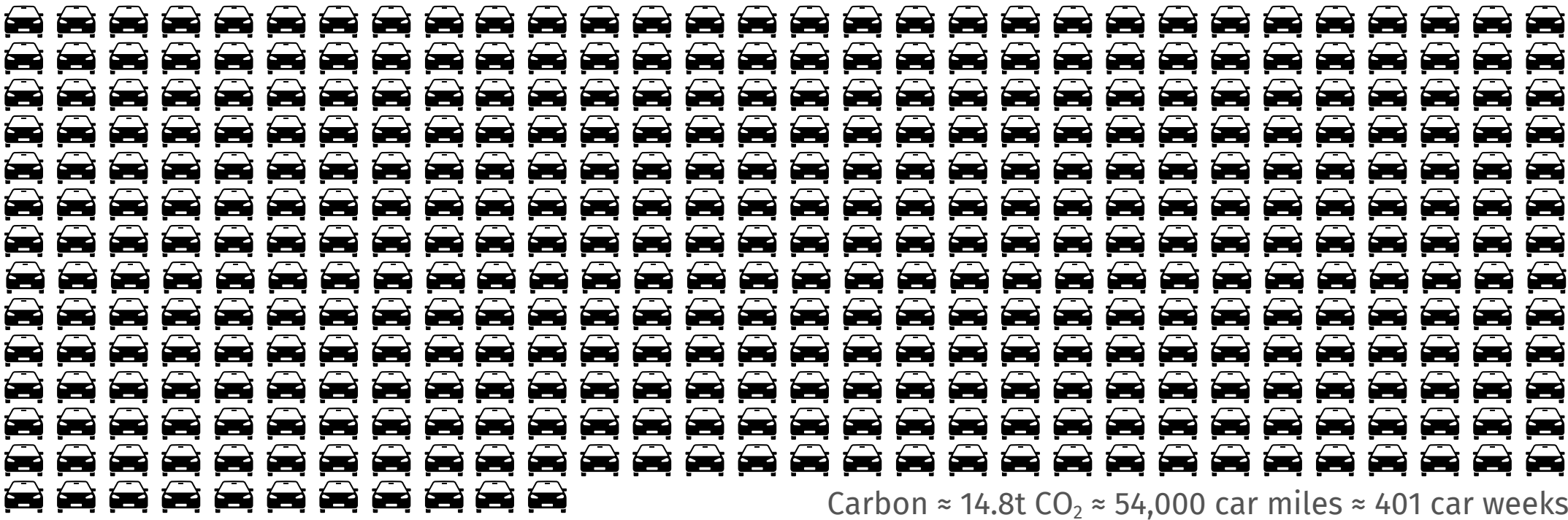
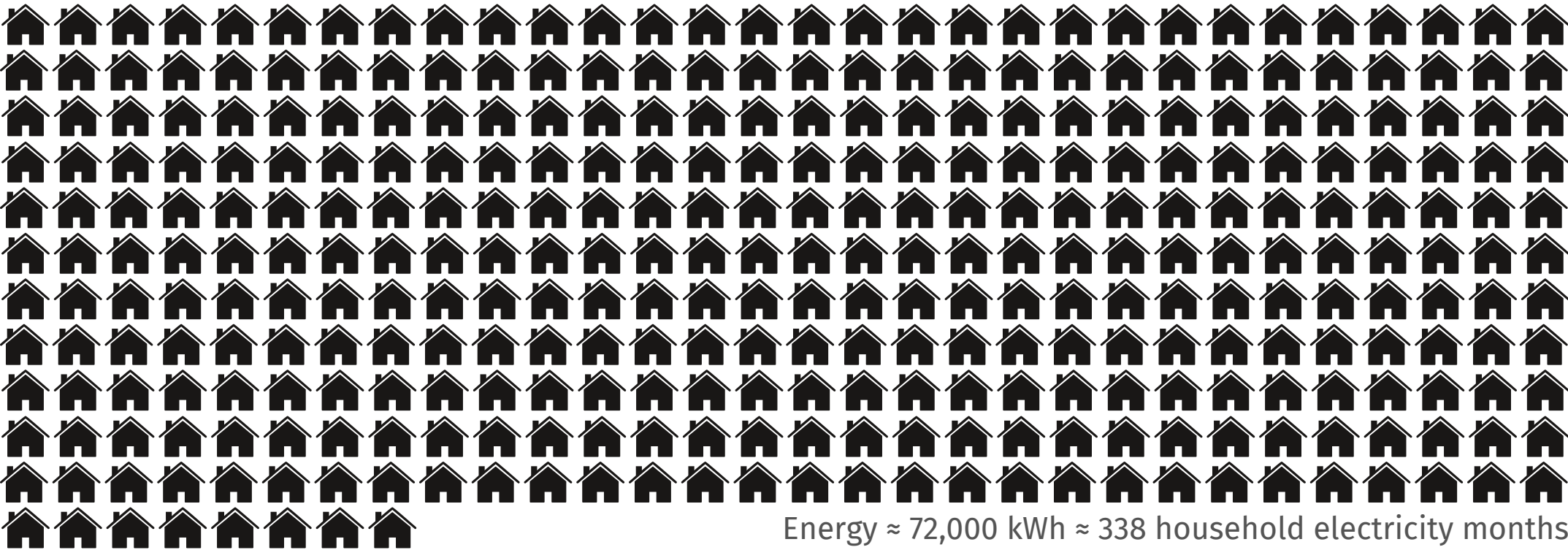
157D vs 159D models:	$\log(z_{\Lambda\text{CDM}})$	$\log(z_{w_0w_a\text{CDM}})$	$\log \text{BF}$	Total computation time
Classical	Unfeasible	Unfeasible	Unfeasible	12 years projected (48 CPUs)
AI-accelerated (ours)	$406689.6^{+0.5}_{-0.3}$	$406687.7^{+0.5}_{-0.3}$	$1.9^{+0.7}_{-0.5}$	8 days (24 GPUs)

Same trained emulator as used previously

(Simulating training data = 1 CPU day | Training = 1 GPU hour | Amortized over all analyses)

Euclid-Rubin-Roman (3x Stage IV survey)-like analysis

Traditional approach



Accelerated approach (ours)



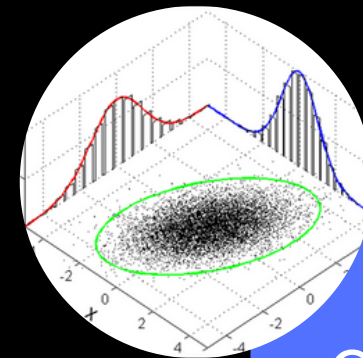
Accelerated scientific inference for physical systems



1st Pillar:
AI
Emulation



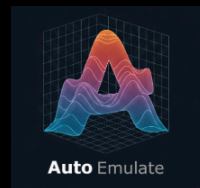
2nd Pillar:
Programming
Frameworks



3rd Pillar:
Gradient-Based
MCMC Sampling



4th Pillar:
Decoupled Model
Comparison



<https://www.autoemulate.com>



<https://github.com/astro-informatics/harmonic>

Dramatic reductions in compute cost, energy usage and carbon emissions... for **every analysis**.