



AI for science, science for AI



# Machine Learning Assisted Sampling: Applications to Physics

Oct 3-7, 2022

ASSAI semester IA-PhyStat  
Workshop summary

Organizers: Marylou Gabrié (CMAP, École Polytechnique),  
Tony Lelièvre (ENPC, Cermics), Valentin de Bortoli (ENS, Deepmind)

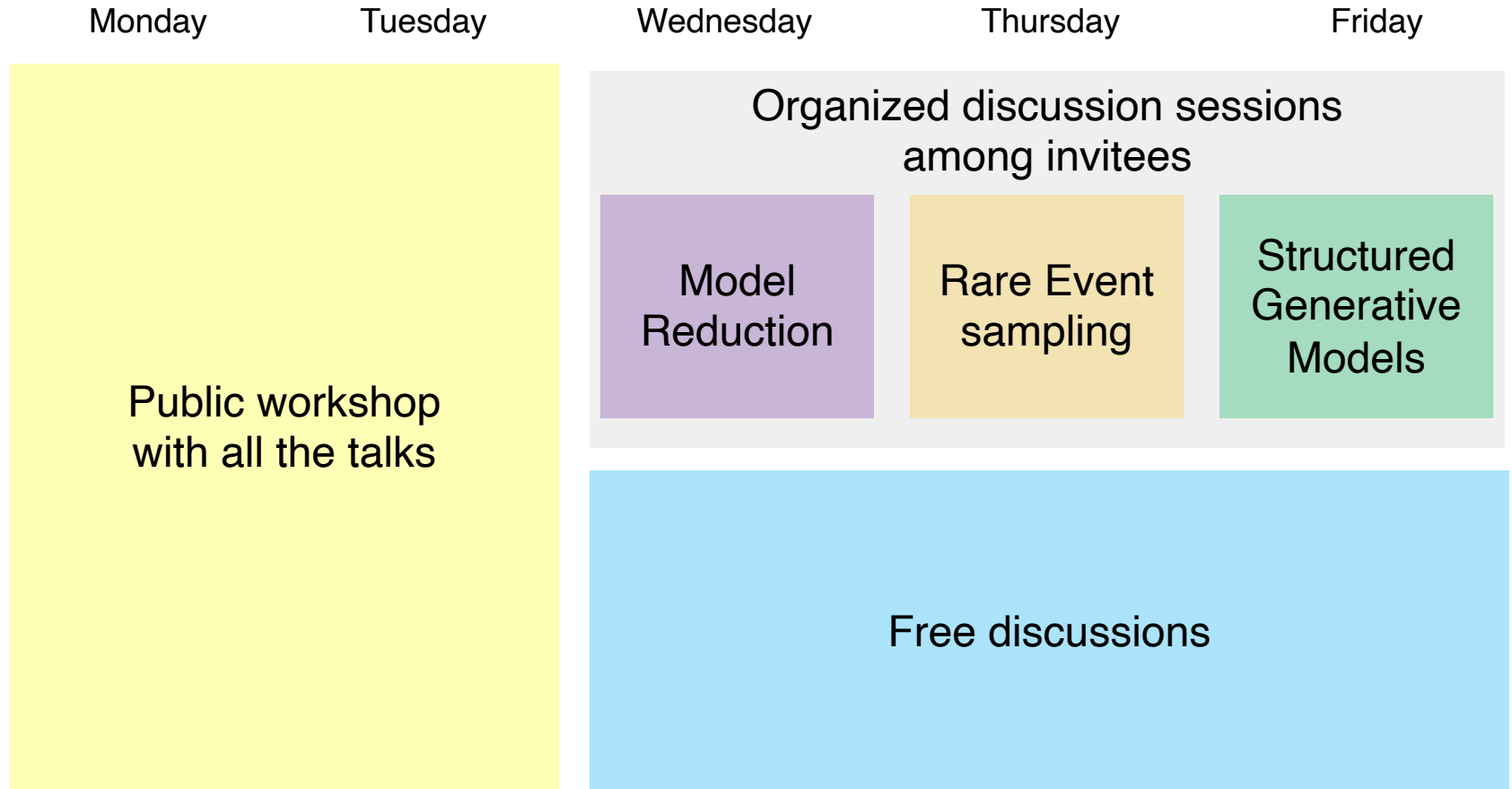
Nov 29, 2023

## **Invited Speakers**

- David Aristoff, Colorado State University
- Peter Bolhuis, University of Amsterdam
- Freddy Bouchet, École Normale Supérieure, Lyon
- Maria K. Cameron, University of Maryland
- Arnaud Doucet, University of Oxford
- Alain Durmus, ENS Paris Saclay & École Polytechnique
- Stéphane Mallat, Ecole Normale Supérieure, Paris
- Pierre Monmarché, Sorbonne Université
- Jutta Rogal, New York University
- Phiala Shanahan, Massachusetts Institute of Technology
- Gabriel Stoltz, École Nationale des Ponts et Chaussées, Paris
- Jonathan Weare, New York University
- Martin Weigt, Sorbonne Université
- Wei Zhang, Zuse Institute Berlin

**Registered participants ~ 80 (50 % online)**

# Format: 2 open days + 3 discussion days

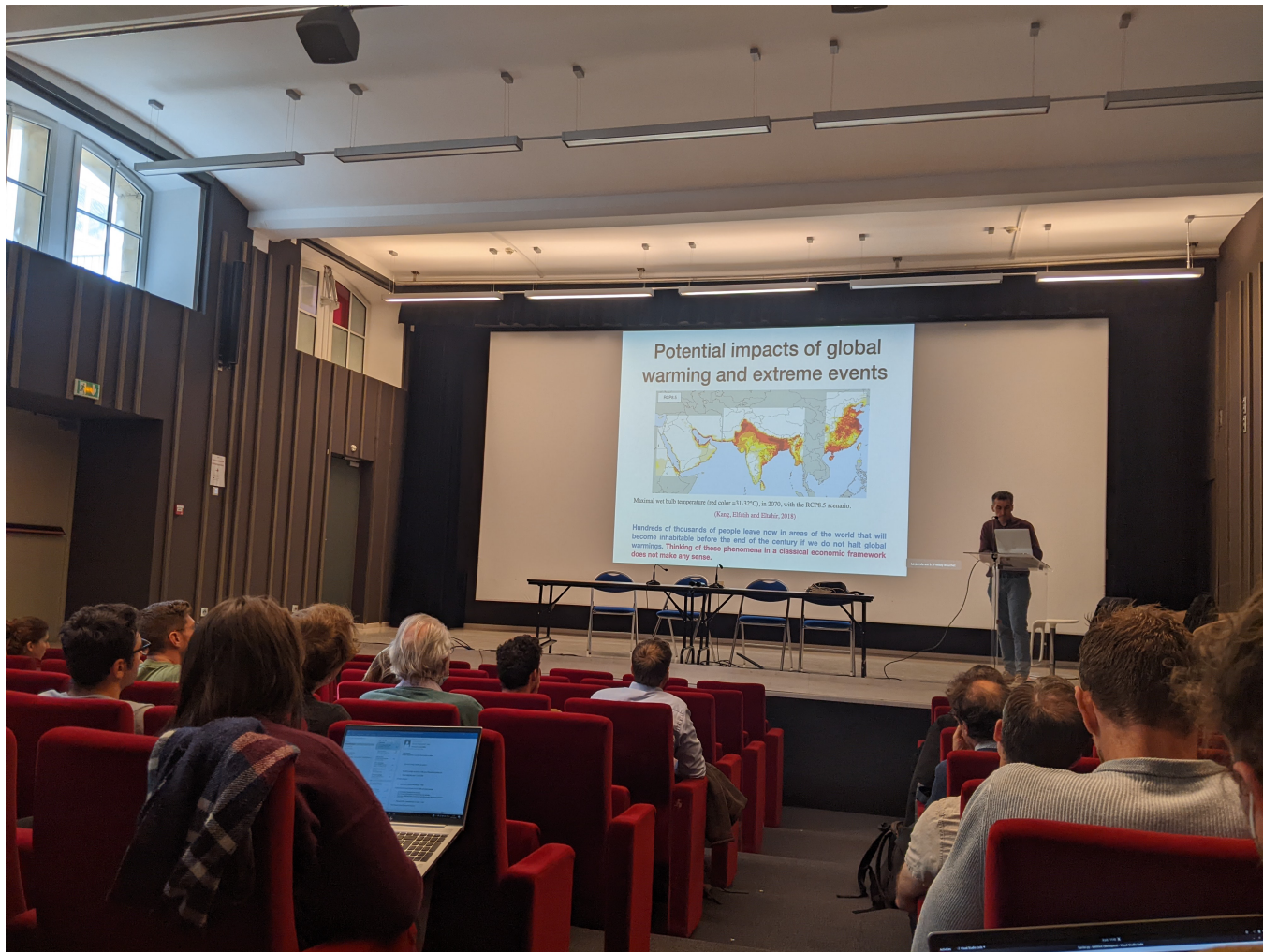


On YouTube!

# Freddy Bouchet's Colloquium:

3

*Probabilistic forecast of extreme heat waves using convolutional neural networks and rare event simulations*



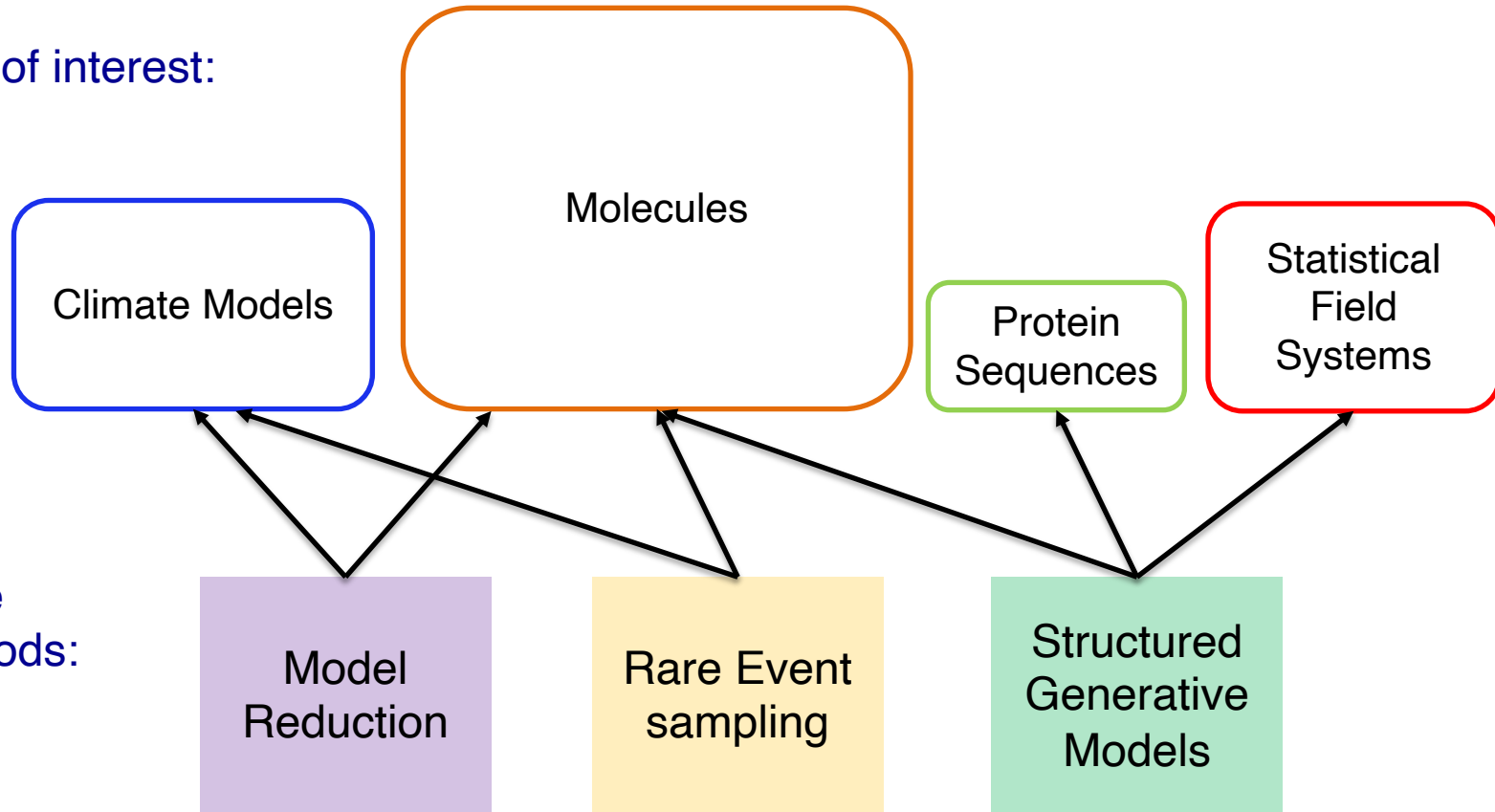
the only picture I took ...

# Main themes & domains of interests of participants

▷ Diverse backgrounds of physicists, applied mathematicians and chemists

▷ All with an interest in sampling

▷ Systems of interest:

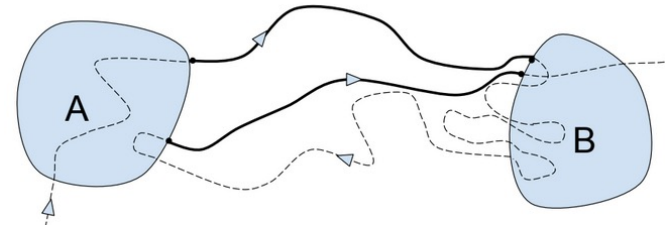


▷ Expertise & methods:

Well established field with experts thinking of opportunities in using machine/deep learning

▷ Given a known dynamic operator and/or scarce observations:

- How to sample rare transitions?
- Identifying reactive channels & maximum likelihood transition paths



- Estimating reaction rates or expected escape times from basins of attractors

*e.g. Understand isomers transformations, characterize whether a storm will become a hurricane*

▷ Important mathematical objects to characterize the dynamics

- Infinitesimal generator

$\mathcal{L}$  = generator of the SDE

- Committed

$q(x) = \mathbb{P}(\tau_A < \tau_B) = \text{committed}$

solution of

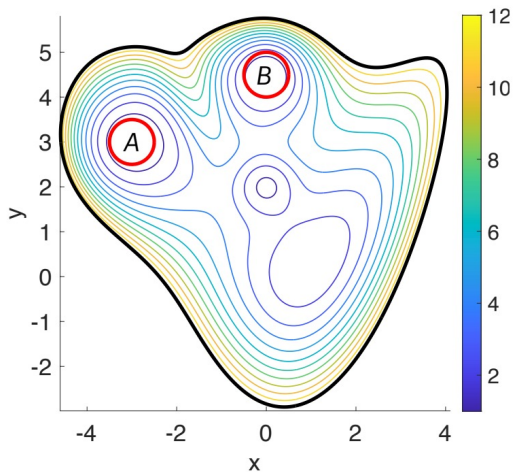
$$\begin{aligned} \mathcal{L}q(x) &= 0, & x &\in (A \cup B)^c \\ q(\partial A) &= 0, & q(\partial B) &= 1 \end{aligned}$$

► Important mathematical objects to characterize the dynamics

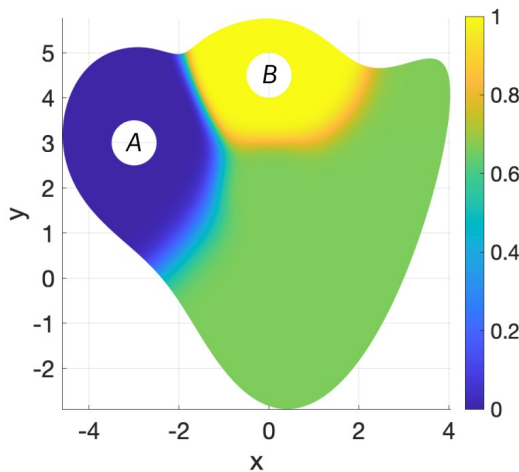
- Infinitesimal generator  $\mathcal{L}$  = generator of the SDE
- Committor  $q(x) = \mathbb{P}(\tau_A < \tau_B)$  = committor

$$\begin{aligned} \mathcal{L}q(x) &= 0, & x \in (A \cup B)^c \\ q(\partial A) &= 0, & q(\partial B) = 1 \end{aligned}$$

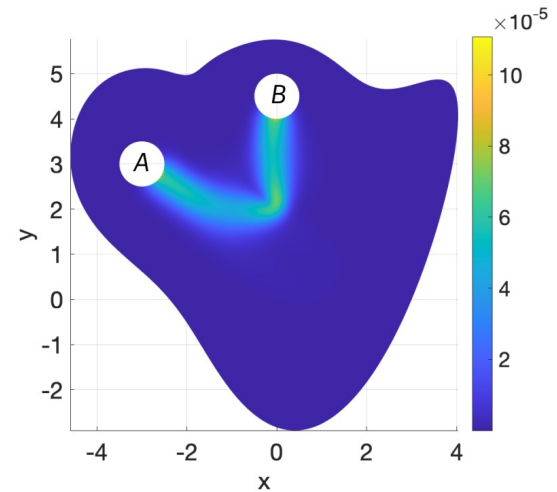
The potential function



The committor



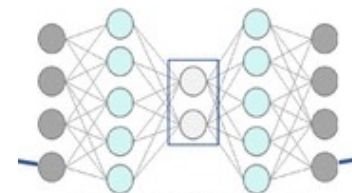
The reactive current



↳ The committor is the optimal “reaction coordinate” to force the transition, but is a high-dimensional function: can it be inferred from data?

## ▷ Dimensionality reductions

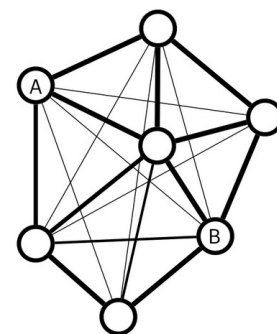
- Use **auto-encoders** to learn non-linear low-d embeddings (Bolhuis, Stoltz)



↳ Which features to include in the reconstruction loss?

- Project the dynamics to discrete states before solving for the committor solving a **linear program** (Cameron, Bouchet)

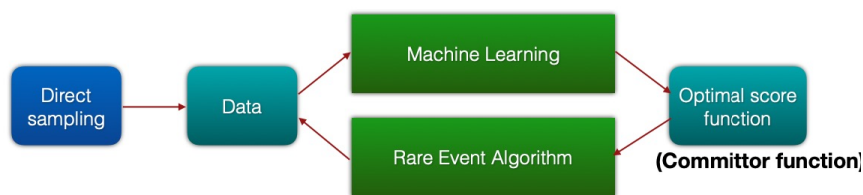
↳ Is the projected dynamics precise enough?



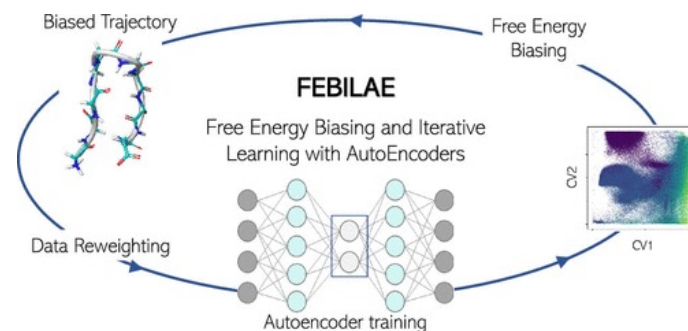
- Identify slow modes associated with metastability using a **feed-forward NN to solve variationally the eigenvalue problem** of the infinitesimal generator (Zhang)

↳ How do NN ansatz compare to usual basis function expansions?

## ▷ Adaptive methods to create data as you go (Bouchet, Stoltz, Gabrié, Weare)



[Credit F. Bouchet]



[Belkacemi et al JCTC 2021]



# Generative models for protein sequence data & statistical field systems

Experts from different fields that have found applications of generative probabilistic models

- ▷ Perform density estimation with a structured model and resample data
  - Predict mutational effects, forecast protein evolutions, design proteins (Weigt)
  - Avoid critical slowing down at resampling time (Mallat)
  
- ▷ Sample from a distribution known up to a normalization constant leveraging a surrogate generative model
  - Bayesian posterior (Doucet, Gabrié)
  - Boltzmann distributions (Shanahan, Gabrié)

# Energy Based and Autoregressive Models for Protein Modelling

▷ Motivations for inferring a sequence landscape

**Data** – MSA of protein family

$$\{(a_i^\mu, \dots, a_L^\mu)\}_{\mu=1, \dots, M}$$



**Statistical sequence model**

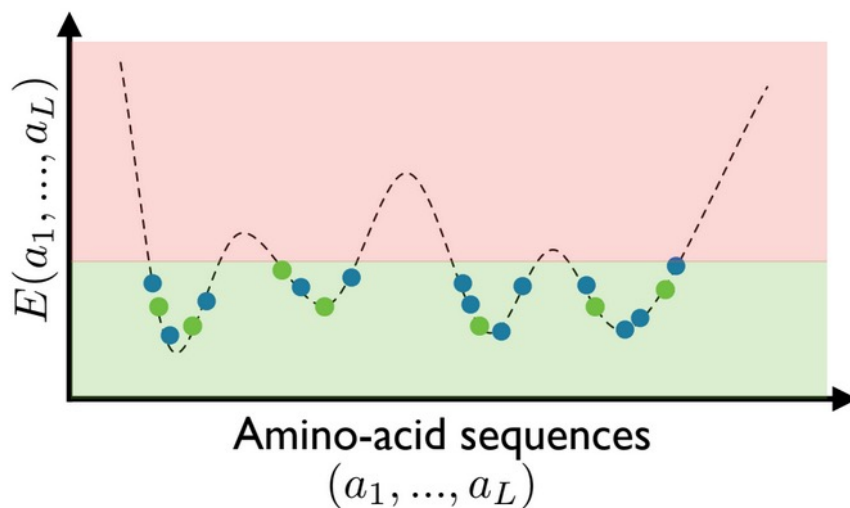
$$P(a_1, \dots, a_L) \sim \exp\{-E(a_1, \dots, a_L)\}$$

“statistical energy”  
“sequence landscape”



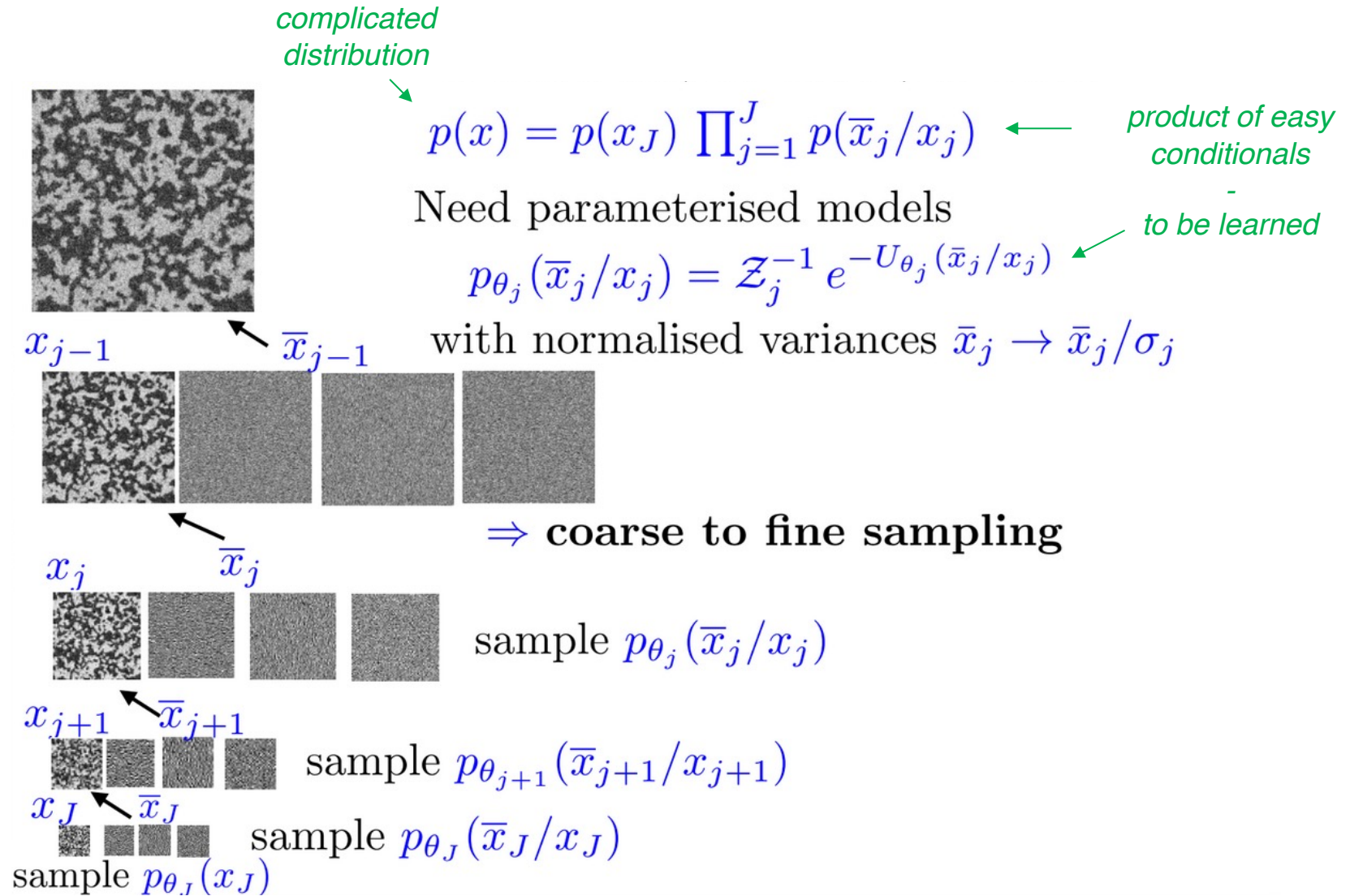
**Sample** – MSA of artificial sequences

$$\{(b_1^\nu, \dots, b_L^\nu)\}_{\nu=1, \dots, N}$$



**Application**

- statistical sequence design



**Log-concave** conditionals allow fast sampling of learned factorization

# Generative models in for protein sequence data & statistical field systems

## ▷ Perform density estimation with a structured model and resample data

- Predict mutational effects, forecast protein evolutions, design proteins (Weigt)
- Avoid critical slowing down at resampling time (Mallat)

↳ Quality of the retrieved model hard to measure

## ▷ Sample from a distribution known up to a normalization constant leveraging a surrogate generative model

- Bayesian posterior (Doucet, Gabrié)
- Boltzmann distributions (Shanahan, Gabrié)

↳ Generative models incorporated in algorithms with performance guarantees

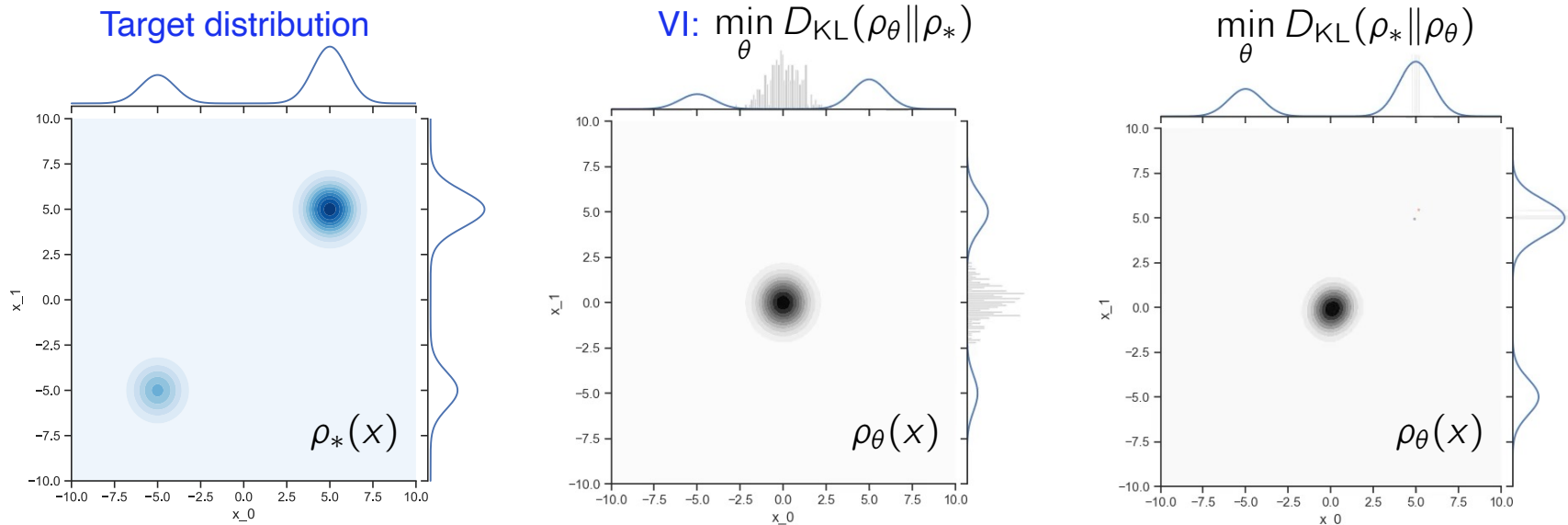
# Assisting sampling with surrogate generative models<sup>12</sup>

No data a priori, only a density of probability  $\rho_*(x)$  (Bayesian posterior, Boltzmann distribution)

▷ **Architecture strategies:** Design generative models to incorporate known symmetries to ease the learning of a surrogate  $\rho_\theta \approx \rho_*$  (e.g. Lattice QCD gauge invariances)

▷ **Training strategies:**

- Variational inference (VI)
- Adaptive training to create data as you go



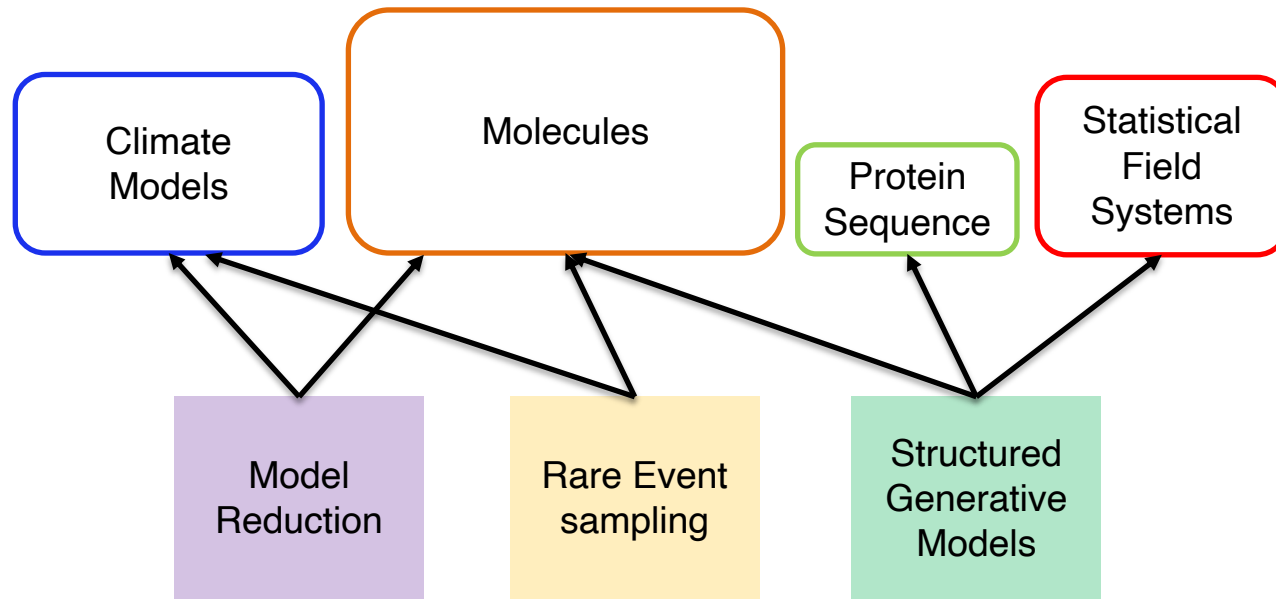
MCMC convergence guarantees!

Adaptive MCMC:

$$\min_{\theta} D_{\text{KL}}(\rho_* \| \rho_\theta)$$

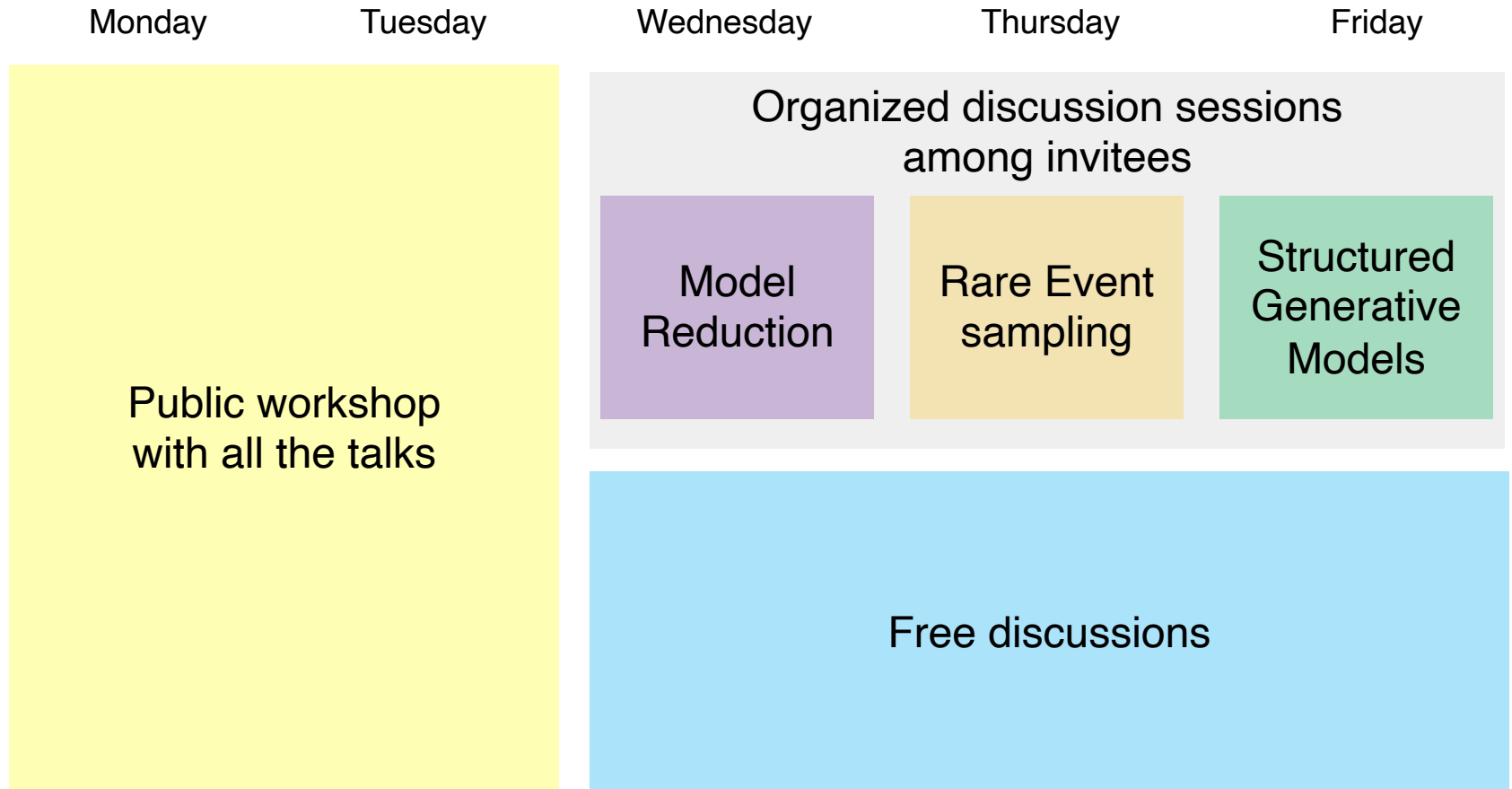
↳ Learning a well covering generative model requires minimum knowledge of modes before-hand

# Conclusion: how/where has machine learning helped?



- ▷ Autoencoders to learn non-linear dimensionality reductions
- ▷ Variational principles to solve eigenvalue problems & partial differential equations
- ▷ Normalized generative models to accelerate sampling
- ▷ Structured generative models to extract/exploit structure from data

# Format: 2 open days + 3 discussion days



On YouTube!



AI for science, science for AI



# Machine Learning Assisted Sampling: Applications to Physics

Oct 3-7, 2022

ASSAI semester IA-PhyStat  
Workshop summary

Organizers: Marylou Gabrié (CMAP, École Polytechnique),  
Tony Lelièvre (ENPC, Cermics), Valentin de Bortoli (ENS, Deepmind)

Nov 29, 2023