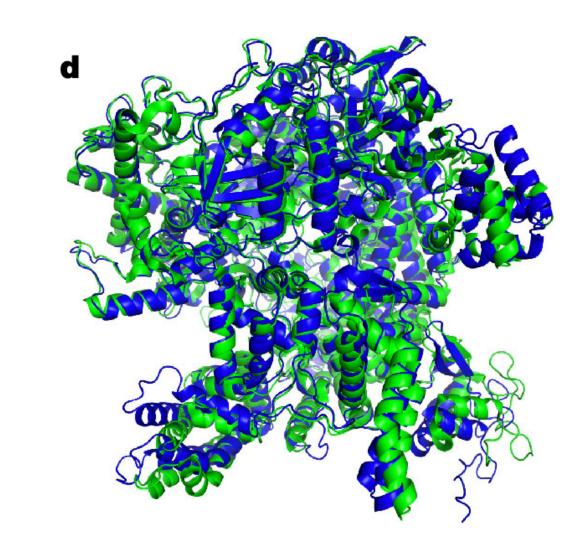
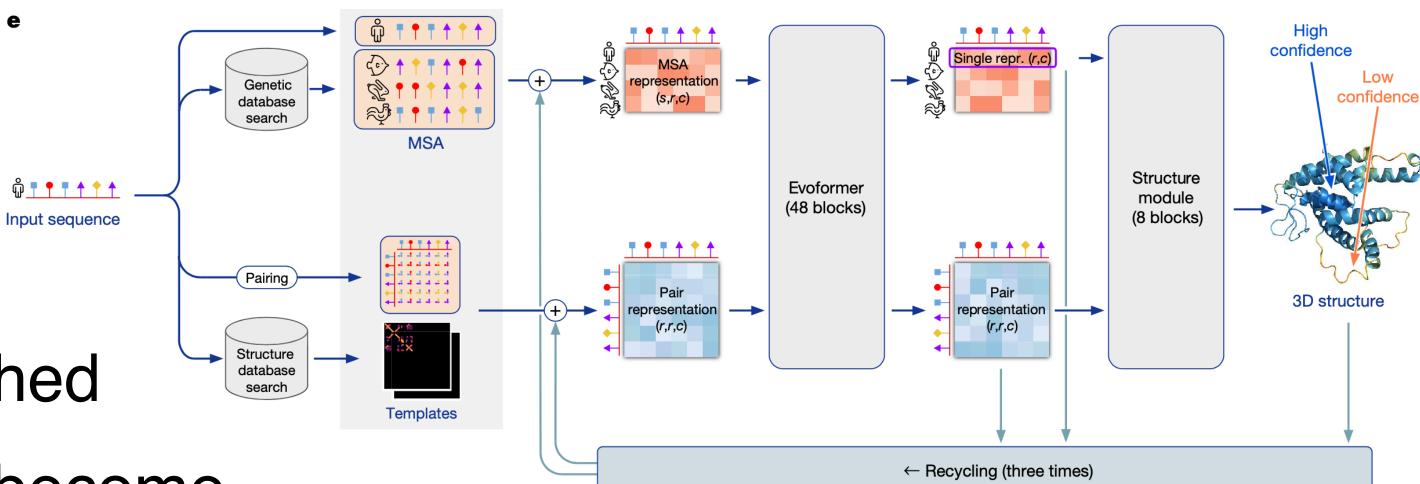
Mathematical Uncertainty Mitigation in Machine Learning (for Science)

AIM NWO Grant Pitch

Motivation

- The success of huge Al models like AlphaFold or ChatGPT, is mainly due to:
 - flexible, scalable models,
 - big engineering effort,
 - huge data sets.
- For such big data sets the attached uncertainties arguably tend to become minimal.





[•] Jumper, J et al. Highly accurate protein structure prediction with AlphaFold. Nature (2021). 2

Problems for typical / scientific data sets

- In contrast, typical or scientific data sets are rather small ($n \lesssim p$ case),
- leading to increased data and model uncertainties and/or complexities,
- requiring sophisticated uncertainty/complexity mitigation strategies,
- necessitating proper mathematical analysis with guarantees!

• We identify **3 major areas** (not necessarily disjoint) for the AIM community to engage in.

Pillar A - Uncertainty / complexity quantification

- Besides the usual model output, also quantify data and model uncertainty/ complexity reliably.
- Possible research directions (examples):
 - statistical/causal hypothesis testing (e.g. e-variables)
 - conformal predictions
 - Gaussian processes
 - statistical and singular learning theory
 - stochastic optimization

Pillar B - Uncertainty / complexity reduction via domain knowledge

- Use domain or expert knowledge to reduce the model complexity and/or data uncertainty.
- Domain/expert knowledge can be conveyed through:
 - models,
 - implicit biases (graph or geometric structures, symmetries, etc.),
 - simulators,
 - differential equations (ordinary/partial/stochastic).
- Possible research directions (examples):
 - geometric and topological machine learning
 - physics-informed machine learning, operator learning
 - simulation-based inference and optimization
 - causal, dynamical and generative modelling

Pillar C - Robustification against uncertainty / complexity

- Identify sources of uncertainty and make the machine learning model output more robust against uncertainty.
- Several approaches and research directions (examples):
 - dimensionality reduction (e.g. PCA, ICA),
 - representation learning (deep, structured, causal),
 - denoising techniques,
 - other data-driven robustifcation techniques (e.g. order statistics, etc.)

Applications in science and society?

- Optional, but desirable are application areas in science and society:
 - medical imaging
 - weather modelling
 - drug discovery
 - molecule and material sciences
 - liquid chromatography
 - astro-physics
 - gene expression data
 - animal movement prediction and tracking
 - Al fairness and safety
 - social networks

The key points

- clear narrative:
 - huge AI models $(n \gg p)$ vs. scientific ML models $(n \lesssim p)$
 - (unreliable) engineering vs. mathematical guarantees
 - optional: useful applications in science and society
 - feasibility proven with our joint previous publication record
- broad umbrella project:
 - easily allows AIM community to participate
 - can be narrowed down (if necessary)

Interested? - What to think about

- To which pillar are you mostly connected to? A, B, C or some combination?
- Think about one PhD student engaging in:
 - one proper mathematical project (Def./Lem./Thm.-style)
 - for publication in (more) mathematical machine learning journal,
 - one project with knowledge transfer from mathematics to machine learning, applying mathematical concepts to machine learning (benchmark data sets)
 - for publication in machine learning conference,
 - one project applying developed methods in science or society (real data sets)
 - for publication in some science journal.

Example - 3D medical imaging

- application area: 3D medical imaging
 - domain knowledge: contains SO(3)-symmetries
 - geometric deep learning (pillar B)
 - study SO(3)-representations, propose SO(3)-equivariant convolutional neural network and study corresponding universal approximation theorem.
 - implement the above neural network, apply it to 3D image/video benchmark data sets and record performance.
 - apply the above method to real 3D medical images for specific prediction or segmentation tasks.

Thank you for your attention!

Questions?