



Machine Learning IN PROTO-NOOS:

Rational Compute Allocation in Antibiotic Discovery via Mutual Information & Uncertainty Quantification

Konrad Gorzelańczyk¹, Szymon Krawczyk¹, Florian Hołubowski¹, Maksymilian Korbik¹

¹Faculty of Computing and Telecommunications, Poznan University of Technology, Piotrowo 3a, 60-965 Poznań, Poland

Correspondence: sknwpl@proton.mail



01 MOTIVATION

Multistage computational pipelines in drug discovery are costly because they often pass all compounds through the same stages, regardless of what is already known from earlier calculations. The motivation behind this project was to test whether mutual information and uncertainty quantification can indicate the point at which further computations can be safely stopped without losing the most important hits. Therefore, the goal was to develop an early-exit strategy for the PROTO-NOOS system, enabling faster exploration of larger chemical space under a limited computational budget.

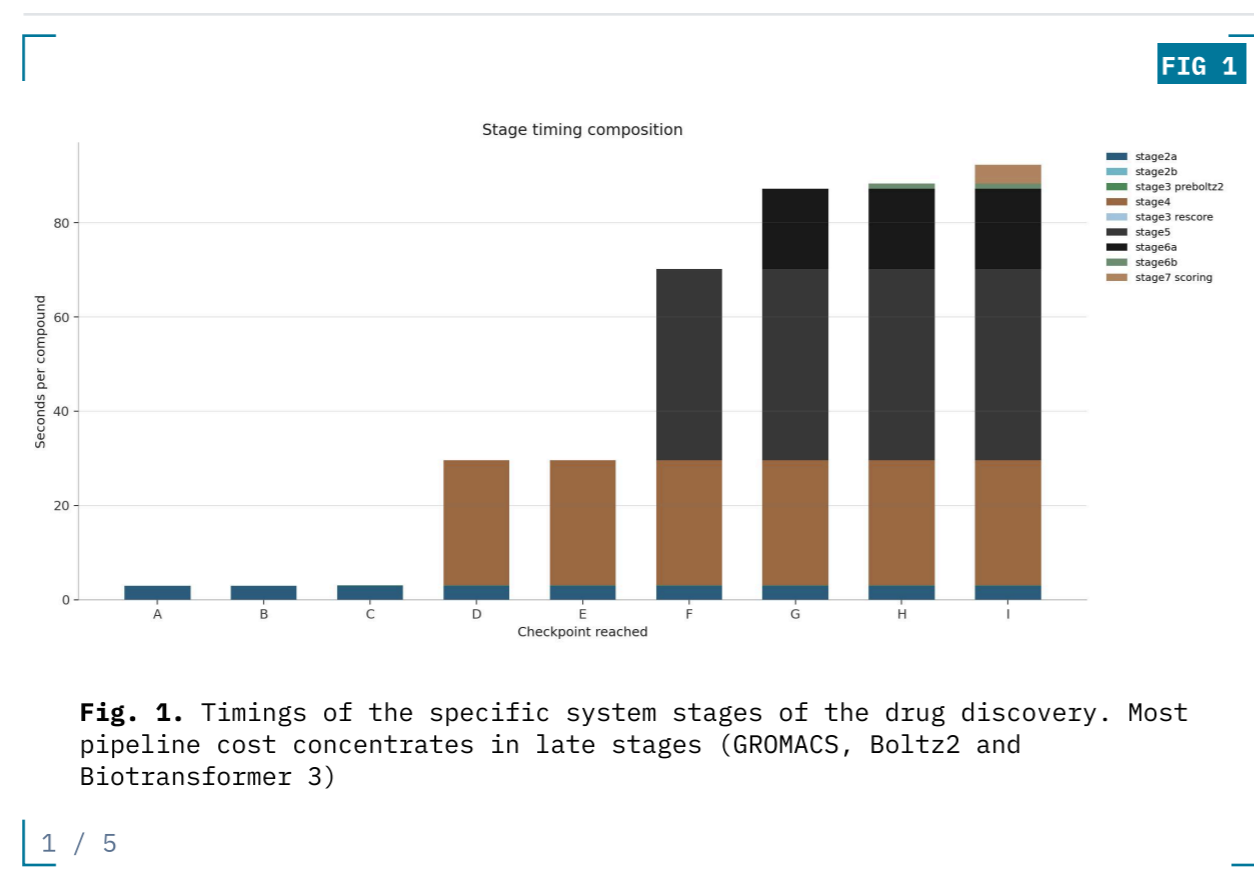
Is knowing when to stop computing as important as knowing what to compute?

779 CURATED INHIB. 2,335 PHENOTYPIC CMPDS 13,531 DE NOVO MOLECULES

02 METHOD

We benchmarked early-exit strategies comparing Single NIG, EOE-5, and EOE-10 surrogate models with uncertainty prediction. Models were evaluated on predefined seeds: 6 for Pareto frontier analysis, 8 for ranking quality, and 10 for UQ diagnostics. Results are reported as mean \pm SD, DDOF = 1. Performance was assessed using hit recovery, compute saved, AUROC, RMSE, Spearman's ρ , distance correlation, and ECE. Mutual information was estimated with MIST-QR and compared against KSG. Result uncertainty was estimated using 1000 bootstrap replicates.

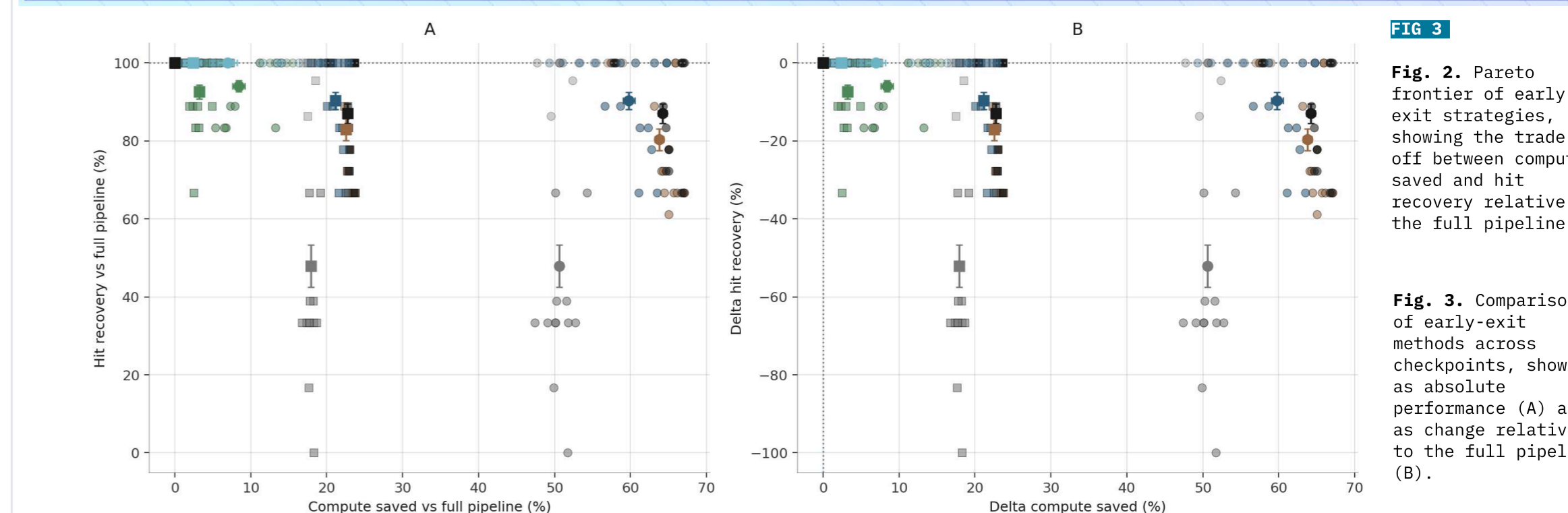
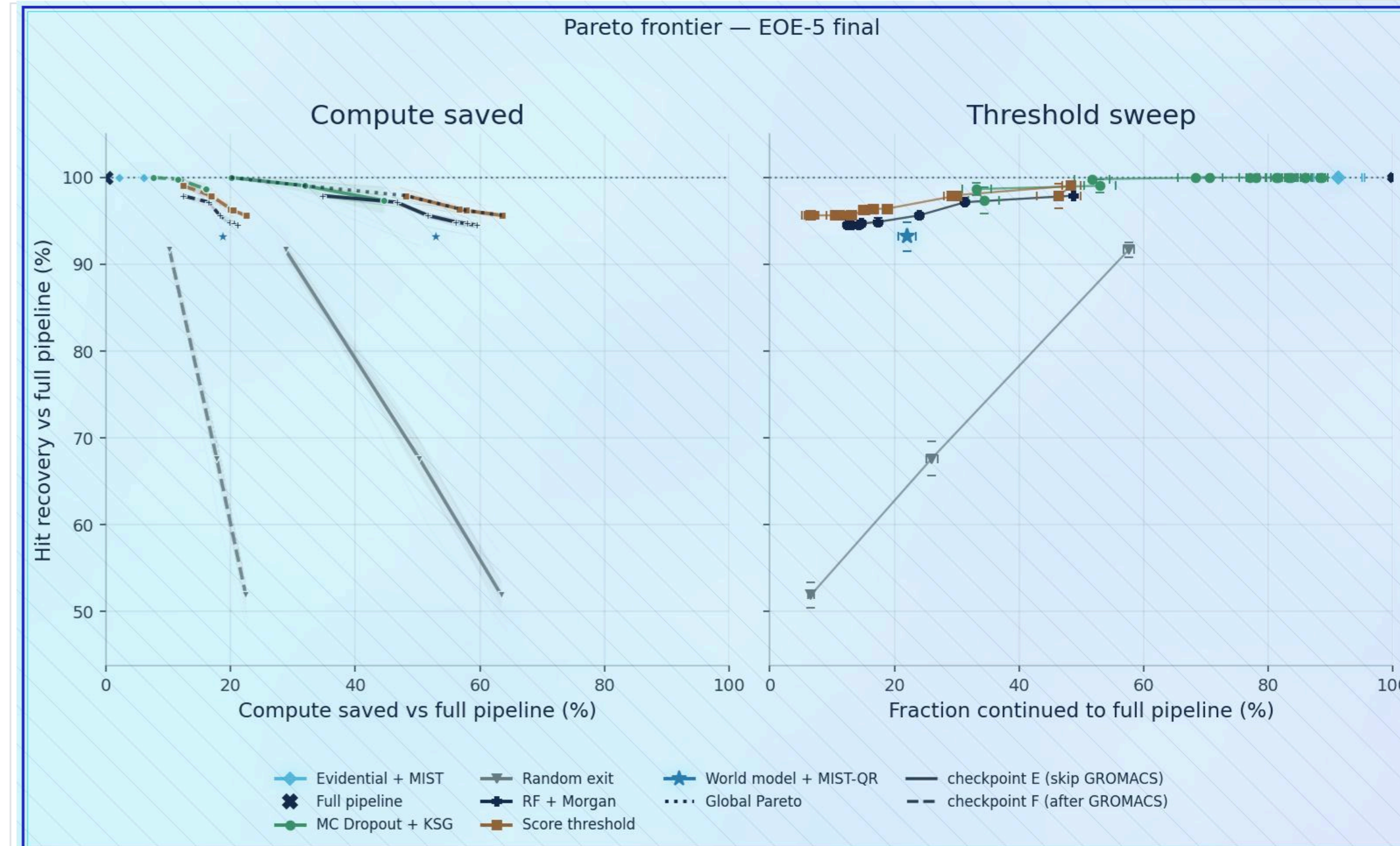
VISUALIZATION 1



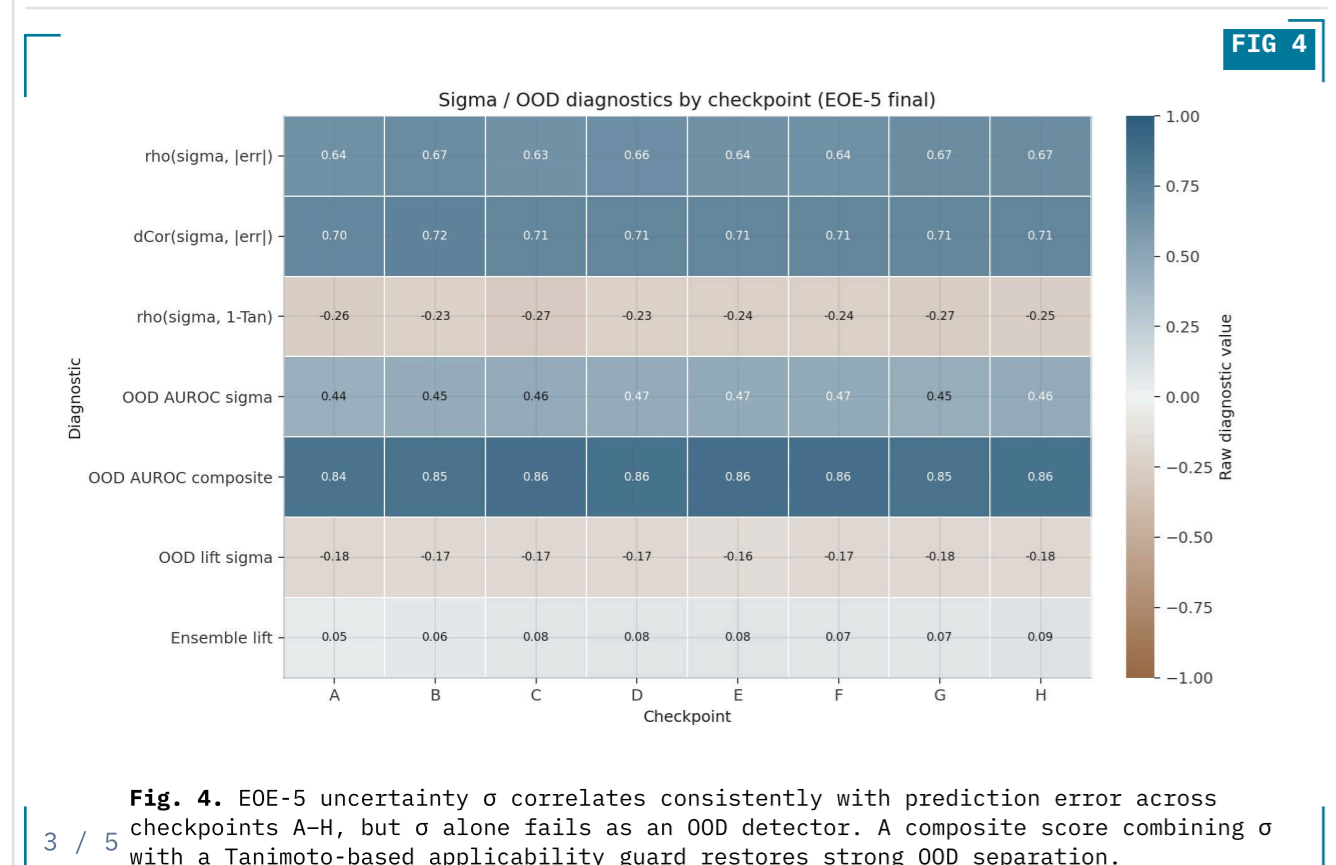
03 MODELS COMPARED

Three model families compared: Single NIG (one forward pass), EoE-5 and EoE-10 (5- and 10-member ensembles of independently trained NIG models). All use UniMol encoder.

VISUALIZATION 2



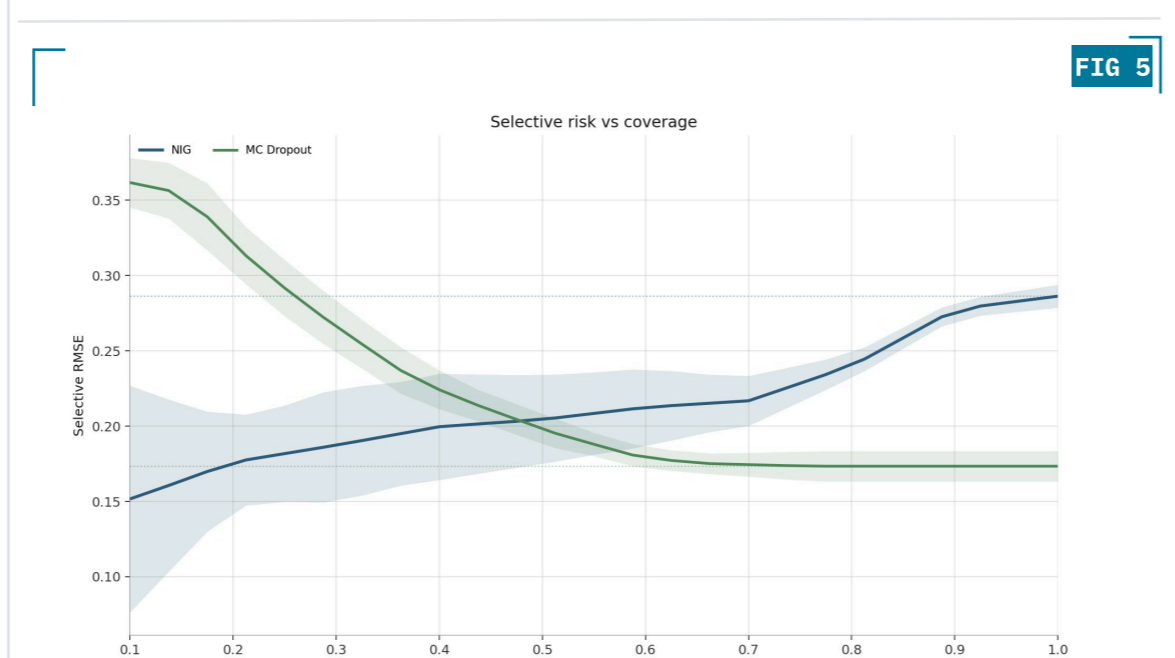
VISUALIZATION 3



04 Pareto Frontier

Evidential + MIST-QR occupies a conservative Pareto point: full hit recovery at ~6% compute saved. More aggressive thresholds (Score threshold, RF + Morgan) save ~58% compute but sacrifice uncertainty reliability in selective prediction.

VISUALIZATION 4



01 TABLE

Method (EoE-5) BEST BOLDED	Compute Saved % (\uparrow)	Hit Recovery % (\uparrow)	AUC-ROC (\uparrow)
Full Pipeline	0.0 \pm 0.0	100.0 \pm 0.0	1.000 \pm 0.000
Random Exit	50.3 \pm 1.5	97.0 \pm 5.5	0.723 \pm 0.029
Score Threshold	57.7 \pm 0.7	100.0\pm0.0	0.945 \pm 0.004
RF + Morgan (Cheap.)	57.9\pm0.1	100.0\pm0.0	0.951 \pm 0.002
MC Dropout + KSG	14.9 \pm 2.6	100.0\pm0.0	0.994 \pm 0.002
World Model + MIST-QR	52.9 \pm 2.4	100.0\pm0.0	0.936 \pm 0.034
Evidential + MIST (Conservative)	6.0 \pm 6.5	100.0\pm0.0	1.000\pm0.001

01 TABLE

Table 1: Best values are bolded. All metrics are averaged over some amount of seeds unless specified otherwise all of the benchmarks were done on the same dataset with same amount of seeds. Averaged over 6 seeds (Pareto), 8 seeds (ranking), 10 seeds (UQ diagnostics) per pre-registered protocol.

02 TABLE

Method BEST BOLDED	Single NIG	EOE-5	EOE-10
Spearman $\rho(\sigma, \epsilon)$ (\uparrow)	0.452\pm0.098	0.409 \pm 0.053	0.444 \pm 0.050
dCor($\sigma, \epsilon $) (\uparrow)	0.514\pm0.076	0.468 \pm 0.046	0.511 \pm 0.042
OOD AUROC σ (\uparrow)	0.576\pm0.043	0.569 \pm 0.045	0.570 \pm 0.037
ECE, 10 adaptive bins (\downarrow)	0.030\pm0.011	0.042 \pm 0.023	0.049 \pm 0.027
MI($\sigma, \epsilon $) KSG (\uparrow)	0.289\pm0.087	0.186 \pm 0.039	0.207 \pm 0.053

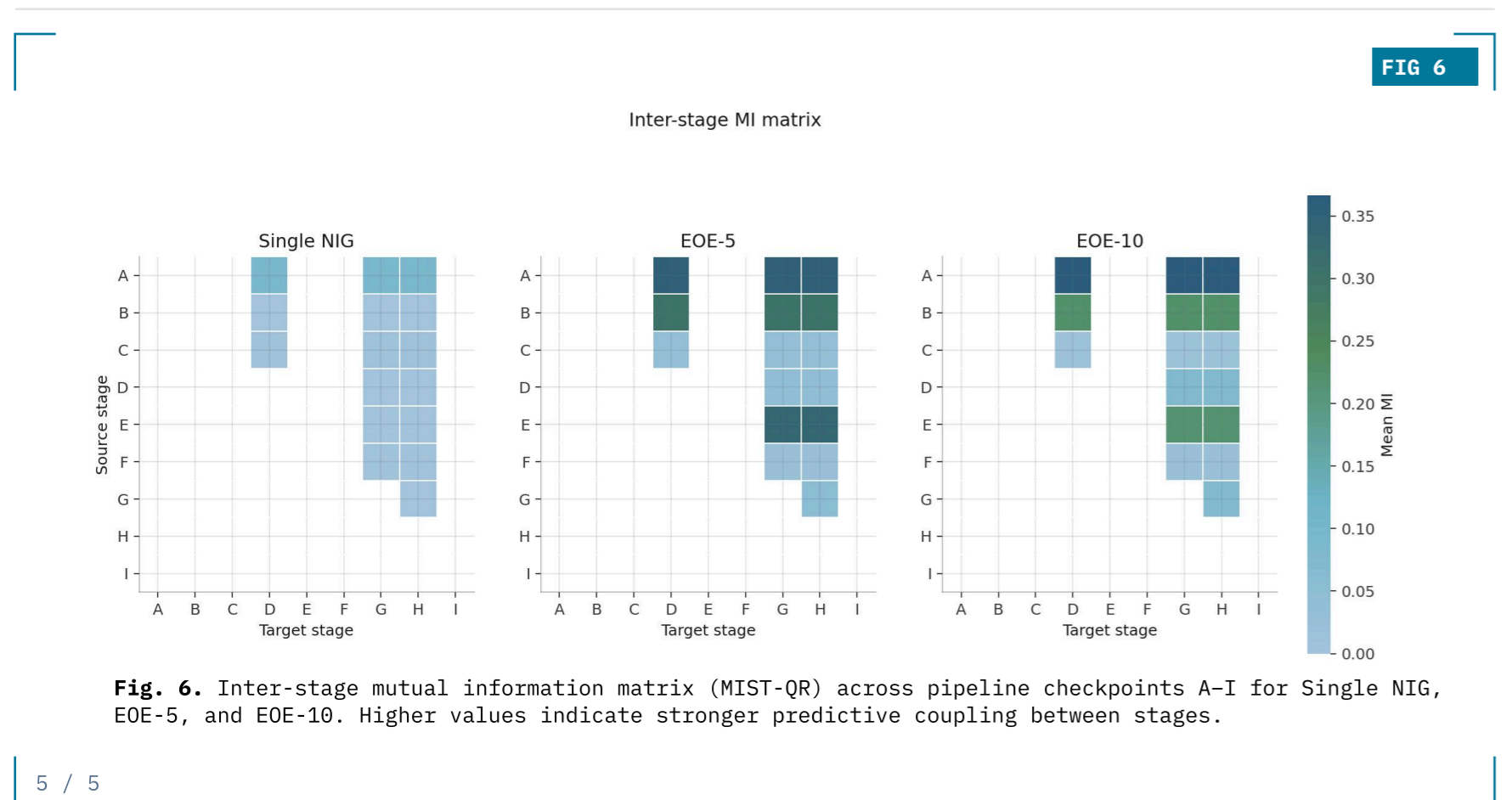
02 TABLE

Table 2: Comparison of metrics from SINGLE NIG, EOE-5, EOE-10. Best are bolded. ECE computed per seed (10 adaptive quantile bins, n=349 per bin), reported as mean \pm SD across 10 seeds. Pooled-prediction ECE shown in Supplementary.

05 OOD / UQ DIAGNOSTICS

NIG uncertainty is useful for ranking prediction risk. However, σ alone gives poor OOD detection. Combining σ with chemical novelty through a Tanimoto guard increases OOD AUROC to 0.84-0.86, showing that model uncertainty and chemical applicability are complementary signals.

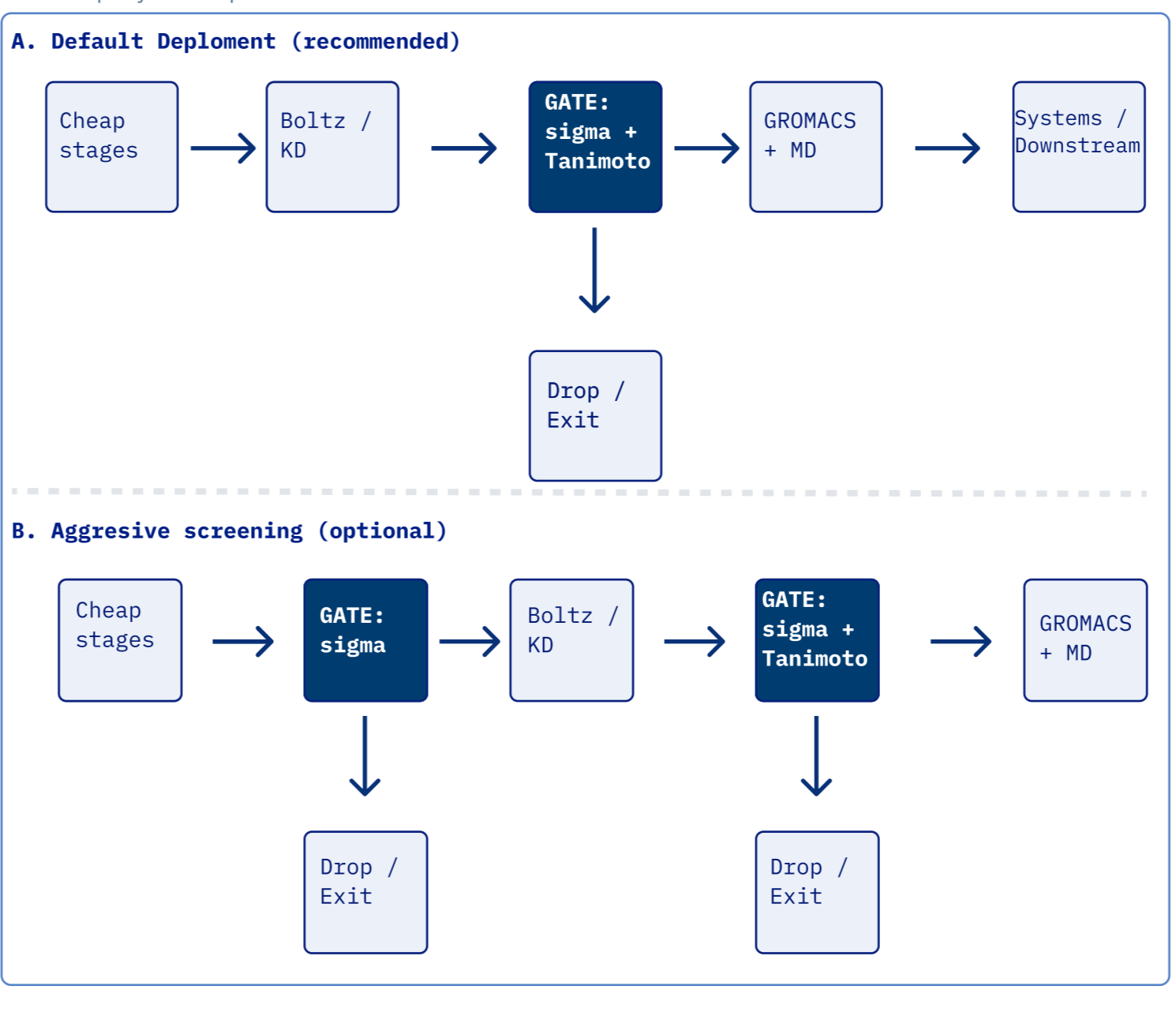
VISUALIZATION 5



06 MI estimates

EOE-5 produces substantially higher inter-stage MI estimates relative to Single NIG, particularly at checkpoints D and E-H. This indicates that ensemble aggregation stabilises mutual information estimation across pipeline stages, rather than simply averaging individual NIG heads. EOE-10 confirms this direction with reduced variance. Single NIG shows near-zero MI for most stage pairs, suggesting that a single forward pass is insufficient to capture inter-stage predictive dependencies reliably.

07 Deployment Options



08 Deployment Steps

- Choose the model: Use Single NIG for production gating. It is cheapest to train and gave the best per-prediction UQ in our benchmark. EOE-5/10 are mainly for diagnostic analysis of inter-stage information flow, not default deployment.
- Place the gate: Default deployment uses one main gate before Boltz/KD and before GROMACS. Optional aggressive screening can add an earlier gate before Boltz/KD.
- Set the threshold: Start with the upper-quantile σ threshold on a held-out set. In our conservative setting, this kept 100% hit recovery with 6% compute saved. Relaxing the threshold can save more compute, but increases risk.
- Add OOD guard: σ alone is insufficient for OOD detection. Use a composite score: $OOD_score = \alpha \cdot 0_nozm + (1-\alpha) \cdot (1 - \max Tanimoto)$ to train, with $\alpha \approx 0.5$. Normalize both terms to $[0,1]$ before combining.
- Validate before deployment: On your held-out set, verify hit recovery, compute saved, and uncertainty calibration before using the gate in production.

09 CONCLUSIONS

Mutual information + uncertainty = principled compute allocation.

- All UQ-based methods preserve 100% hit recovery. The difference between methods is not in hit recovery, but in the trade-off between compute saved and UQ reliability, measured by AUROC and selective risk. This reframes the method ranking beyond the classical "hit recovery vs cost" view.
- Per-prediction UQ: Single NIG is sufficient. It achieves the best Spearman's $\rho(\sigma, |\epsilon|) = 0.452$, $dCor = 0.514$, $ECE = 0.030$, and $MI(\sigma; |\epsilon|) = 0.289$ nats. The ensemble does not improve uncertainty quality for individual molecules. For selective prediction and early-exit gating, Single NIG is the appropriate choice.

- σ alone is not sufficient for OOD detection. AUROC ~ 0.57 across all three models is close to a random baseline. OOD detection requires a composite score combining σ with a chemical applicability guard, such as Tanimoto < 0.30 . Model uncertainty and chemical novelty are complementary signals, not interchangeable ones.

- MIST-QR reveals inter-stage dependencies that are not visible in a single model. Single NIG shows near-zero MI between stages, whereas EOE-5 and EOE-10 provide a meaningful signal. The ensemble not only averages predictions, but also stabilizes MI estimation.

OOD AUROC 0.57 (σ ALONE INSUFFICIENT)

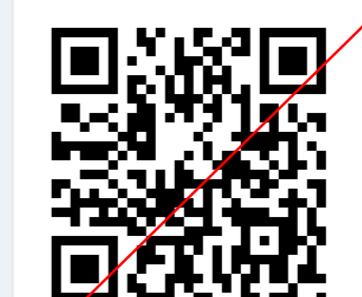
ECE 0.030 (SINGLE NIG)

$I(\sigma; |\epsilon|) = 0.29$ NATS

6% to 58% COMPUTE SAVED

PREPRINT x REPOSITORY

PREPRINT



<https://komputerowe-projektowanie-lekow.github.io/webpage/>

CODE



<https://github.com/>