



AutoML: Replacing Data Scientists?

Marius Lindauer

Introducing myself

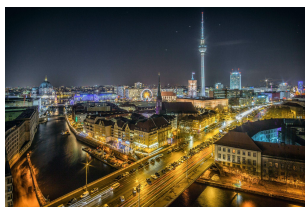
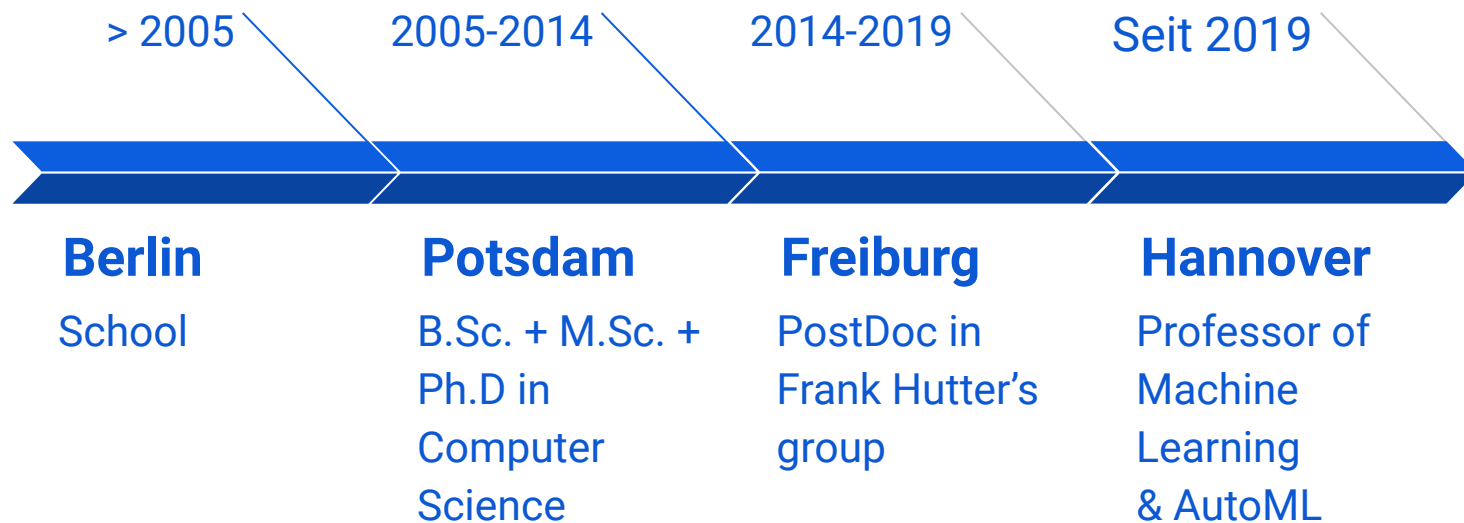


Photo by [Stefan Widua](#) on [Unsplash](#)



Photo by [Chris Revem](#) on [Unsplash](#)



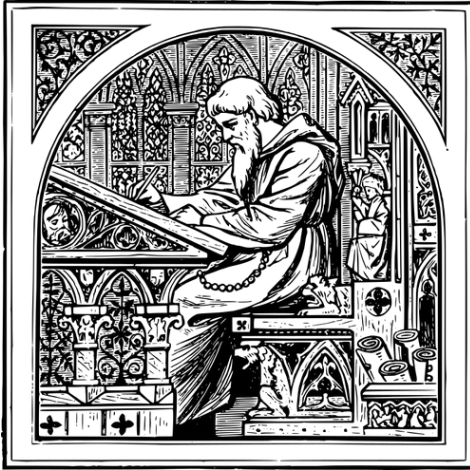
Photo by [Kanan Khasmammadov](#) on [Unsplash](#)



Photo by [op23](#) on [Unsplash](#)

The need of AutoML!?

Rise of Literacy



- Only priests were able to read and write
- People believed that they don't need to read and write
- They went to the holy buildings



Photo by [Anna Hunko](#) on [Unsplash](#)

- Today, everyone can read and write
- No one doubts the benefits of it
- **⇒ Democratization of literacy**

Inspired by [Andrew Ng](#)



Photo by [Max Duzij](#) on [Unsplash](#)

- Only highly educated people can program new AI applications
- Power only with the large IT companies



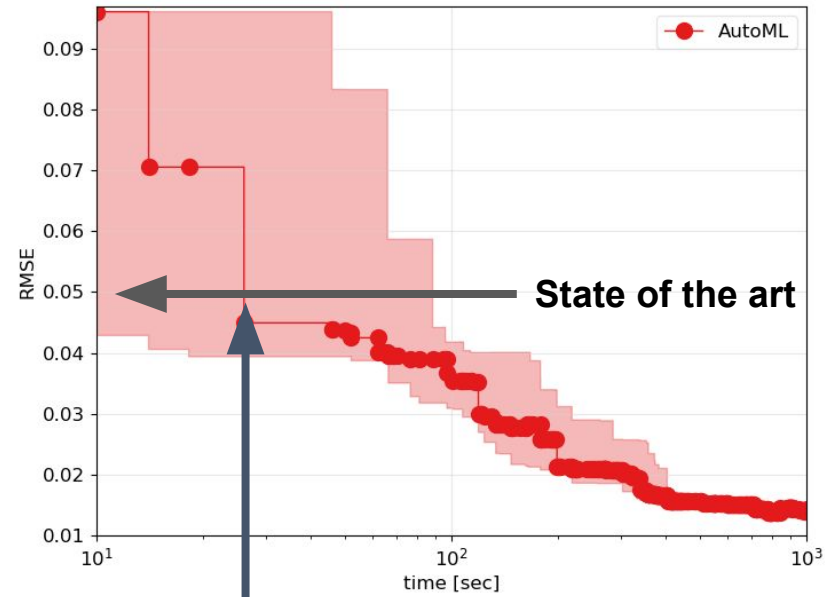
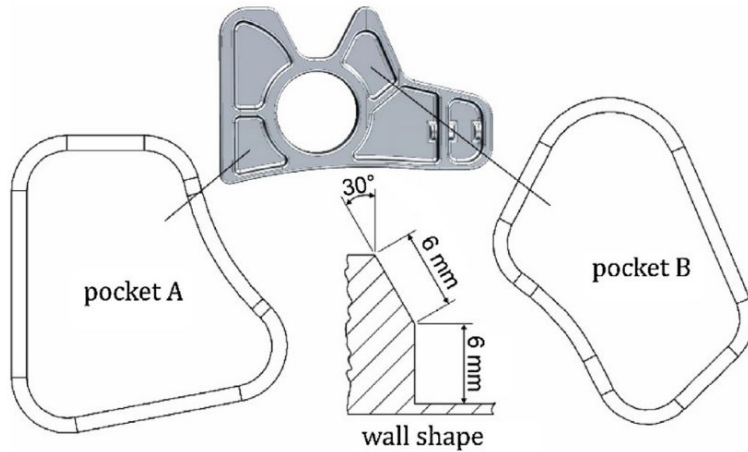
- In an age of limited resources, the need for efficient use gets more important
- **AutoML contributes to AI literacy!**

[\[See also my TEDx Talk\]](#)

A case study with engineers [\[Denkena et al. 2020\]](#)

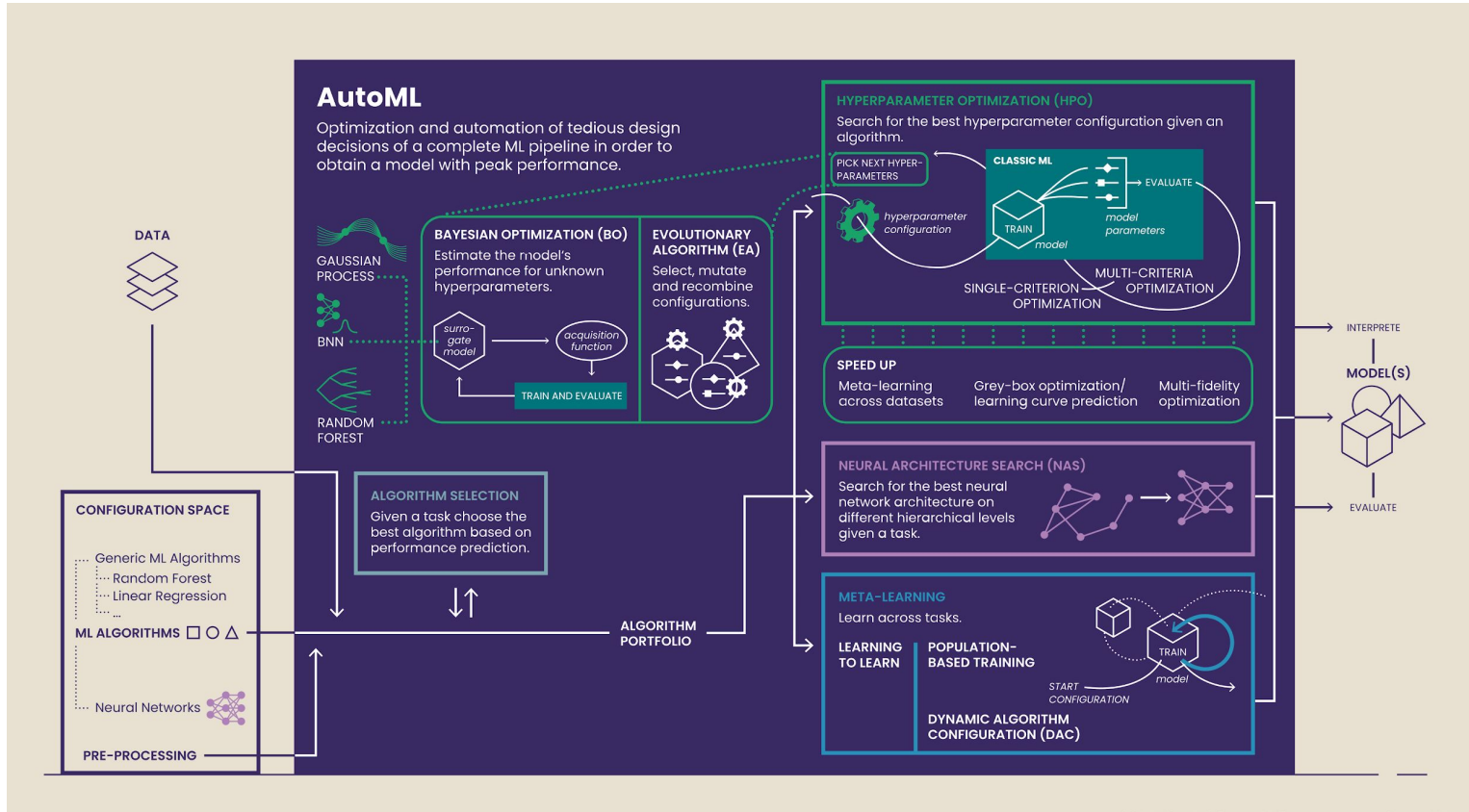


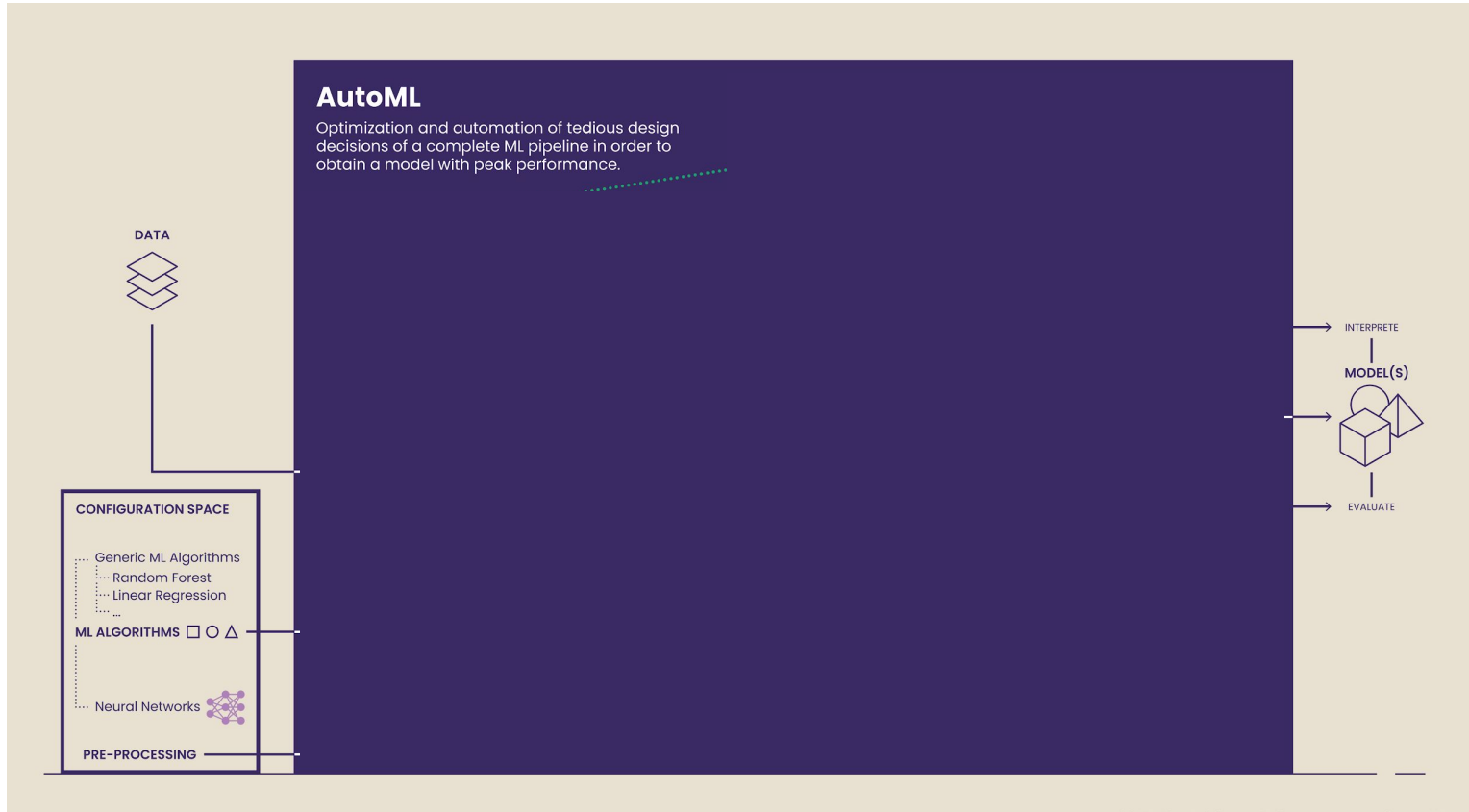
Shape Error Prediction in Milling Processes



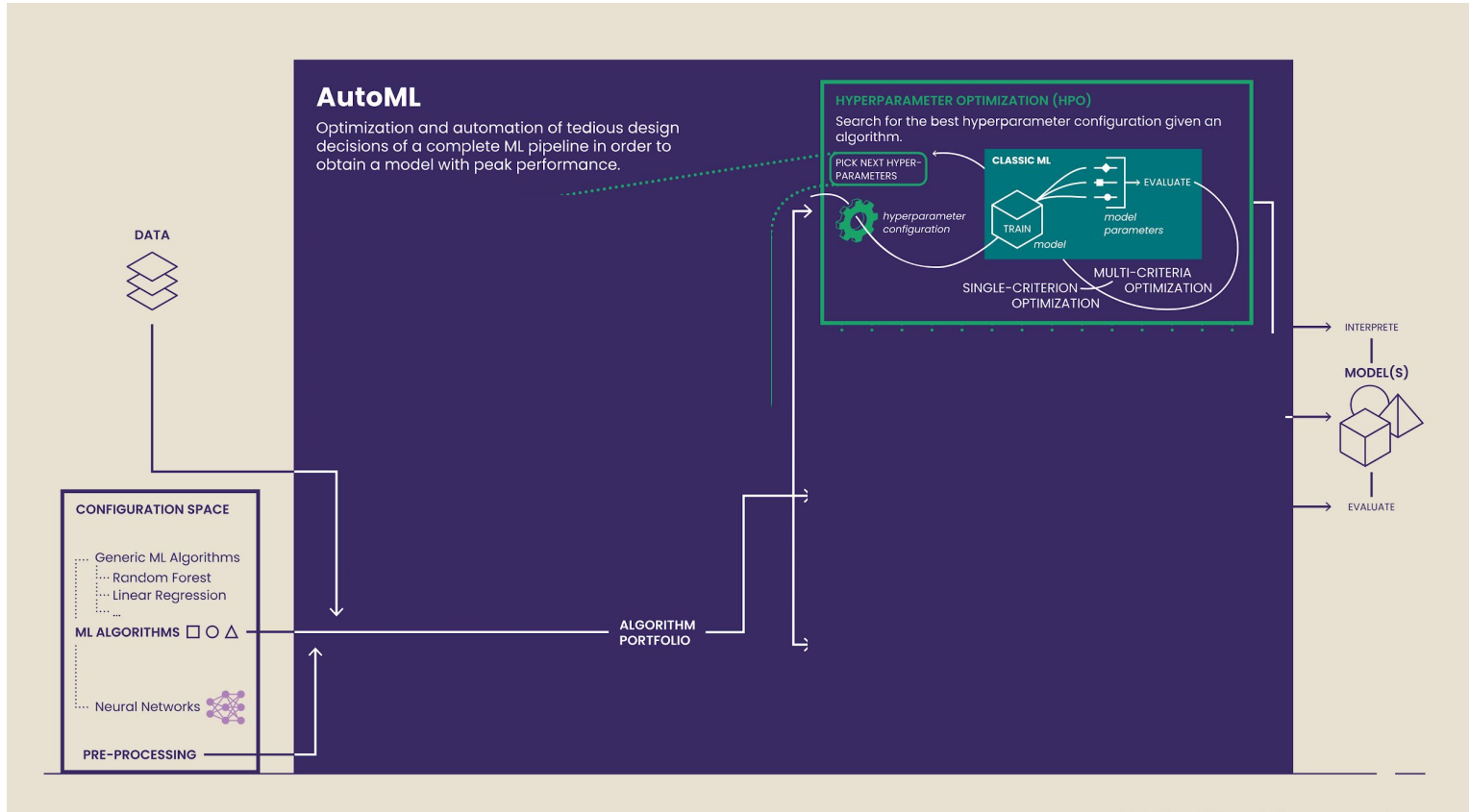
**Better than state of the art
after 27 sec!**

AutoML Landscape





Hyperparameter Optimization



HYPERPARAMETER OPTIMIZATION (HPO)

Search for the best hyperparameter configuration given an algorithm.

PICK NEXT
PARAMETERS

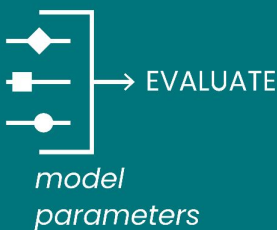


*hyperparameter
configuration*

CLASSIC ML



model



EVALUATE

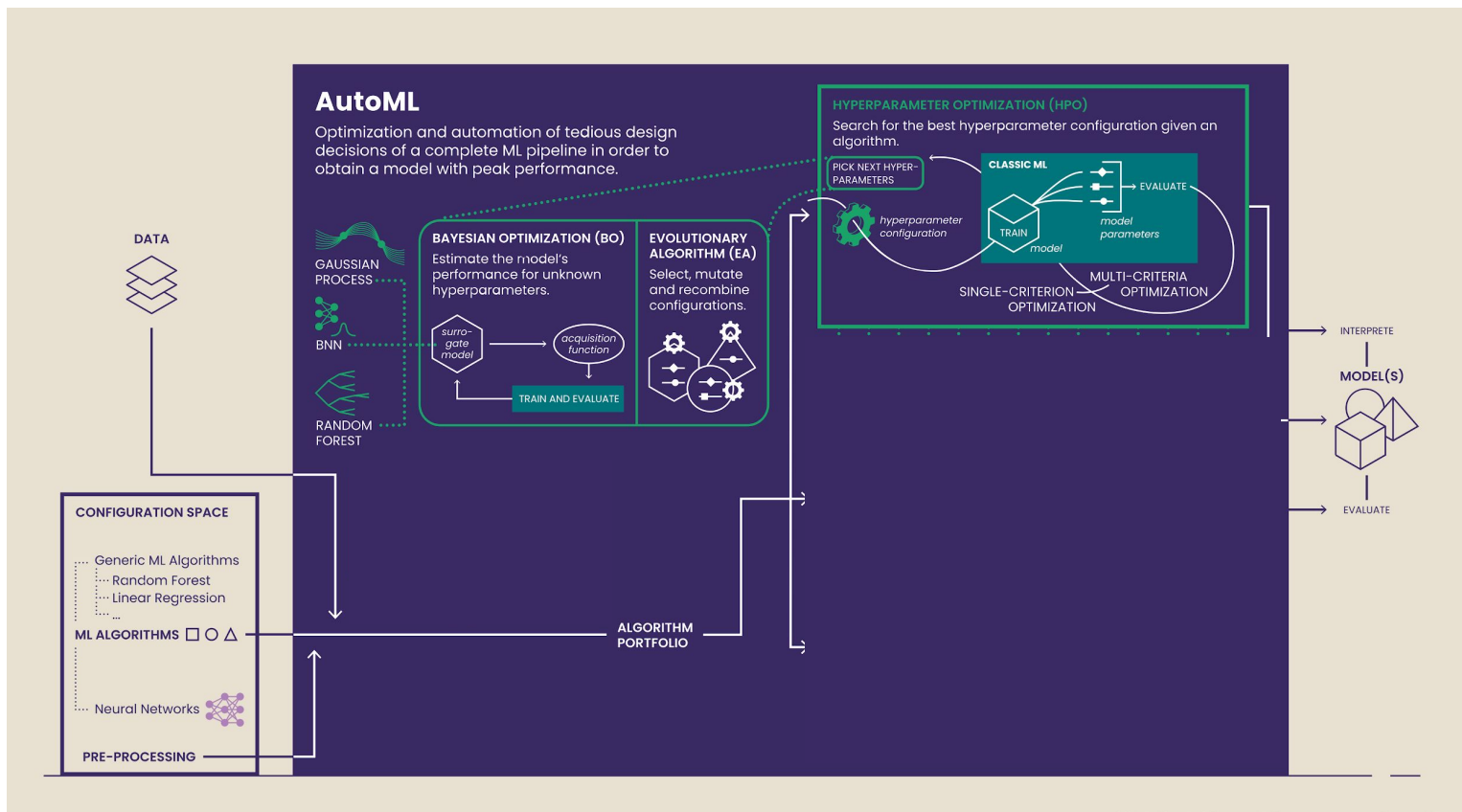
*model
parameters*

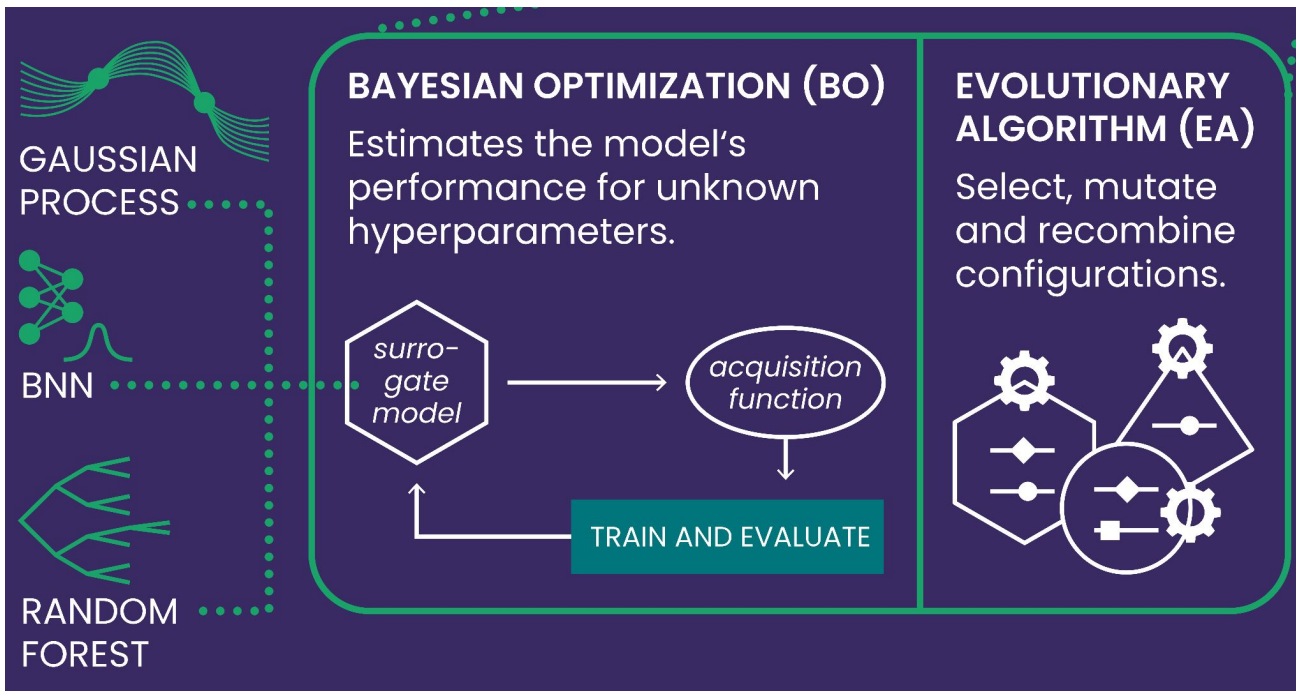
SINGLE-CRITERIA
OPTIMIZATION

MULTI-CRITERIA
OPTIMIZATION

Optimize for

- Accuracy (& co)
- Memory consumption
- Energy consumption
- Inference time
- Training time
- Fairness
- Robustness
- Uncertainty quantification
- ...

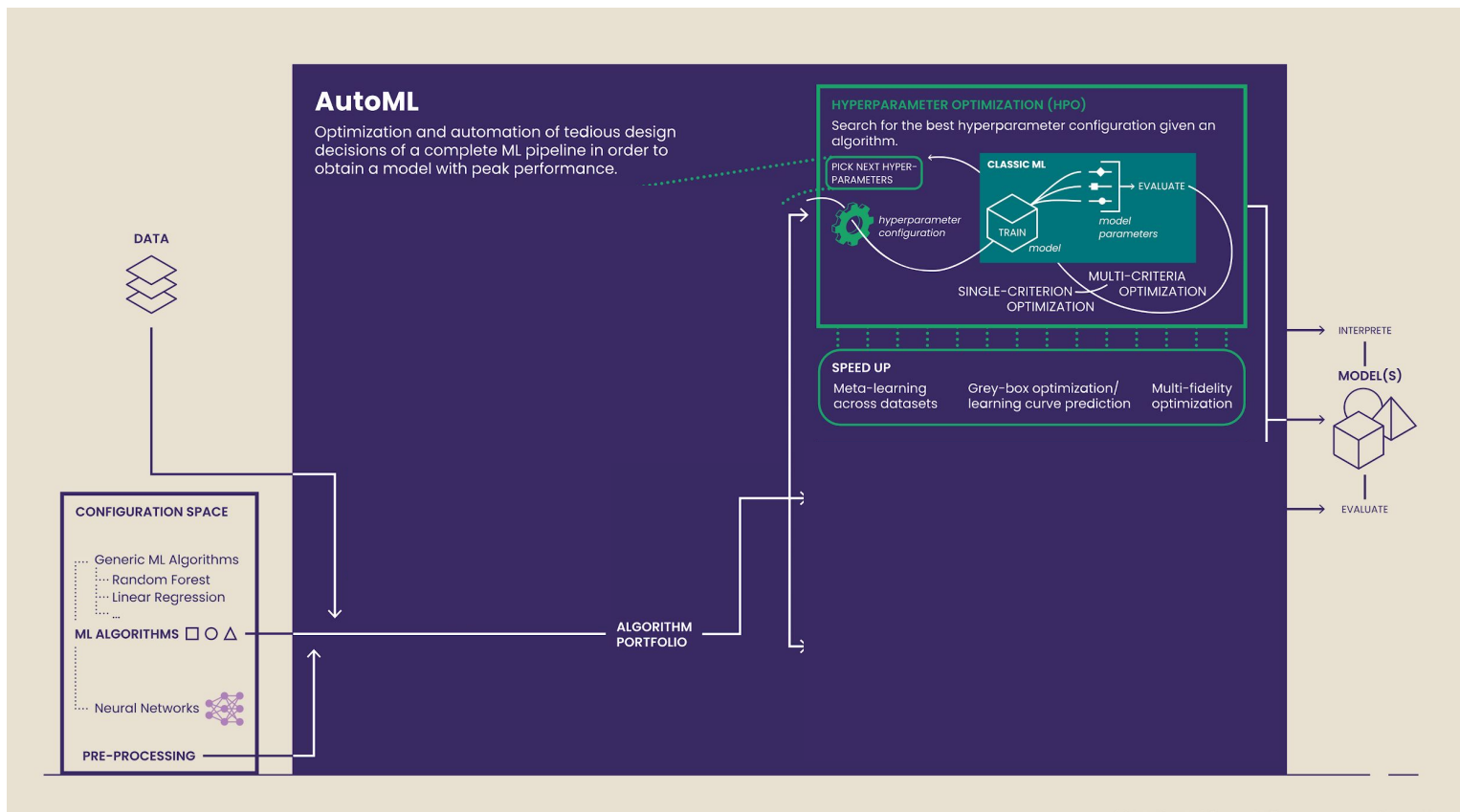




Further alternatives:

- Grid search
- Random search
- Reinforcement Learning
- Planning

Speeding Up

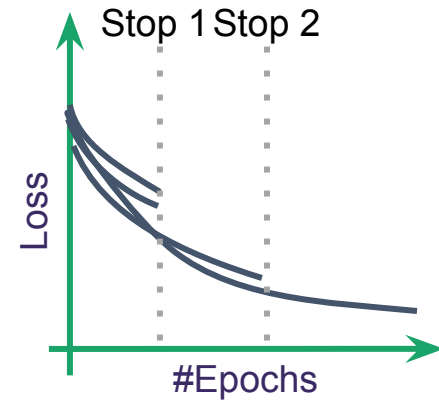
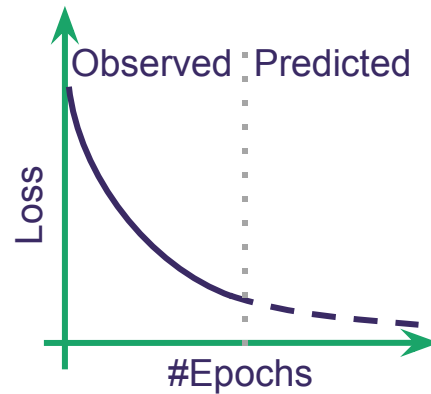
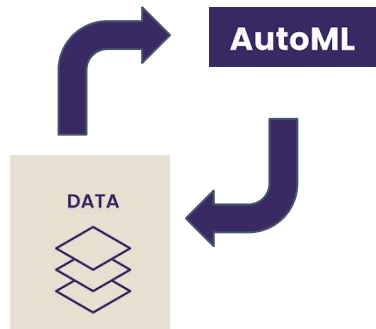


SPEED UP

Meta-learning
across datasets

Grey-box optimisation/
learning curve prediction

Multi-fidelity
optimisation



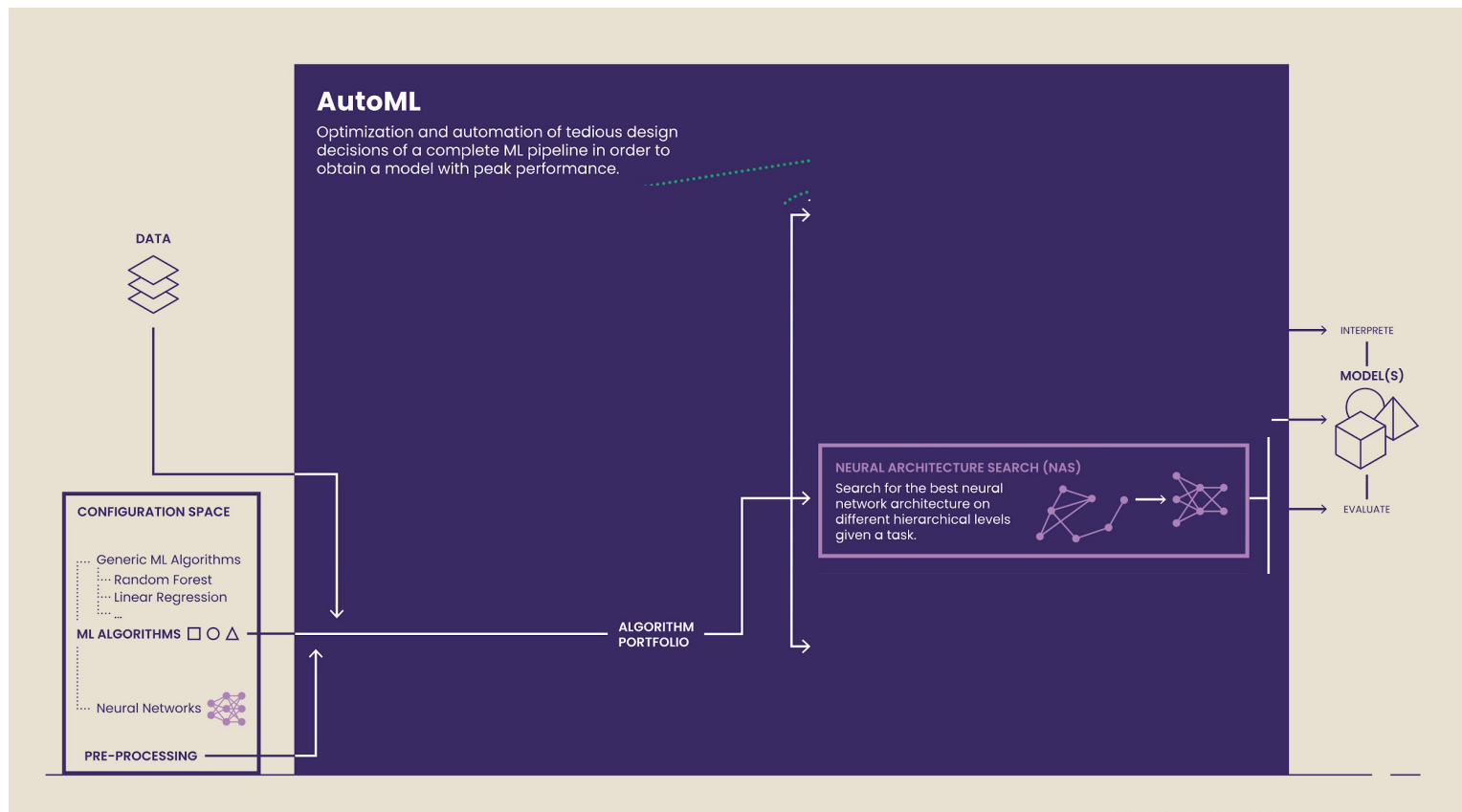
HPO Packages

Package	Complex Hyperparameter Spaces	Multi- Objective	Multi- Fidelity	Instances	CLI	Parallelism
HyperMapper	✓	✓	✗	✗	✗	✗
Optuna	✓	✓	✓	✗	✓	✓
Hyperopt	✓	✗	✗	✗	✓	✓
BoTorch	✗	✓	✓	✗	✗	✓
OpenBox	✓	✓	✗	✗	✗	✓
HpBandSter	✓	✗	✓	✗	✗	✓
SMAC	✓	✓	✓	✓	✓	✓



last update of table in 2021

Neural Architecture Search (NAS)



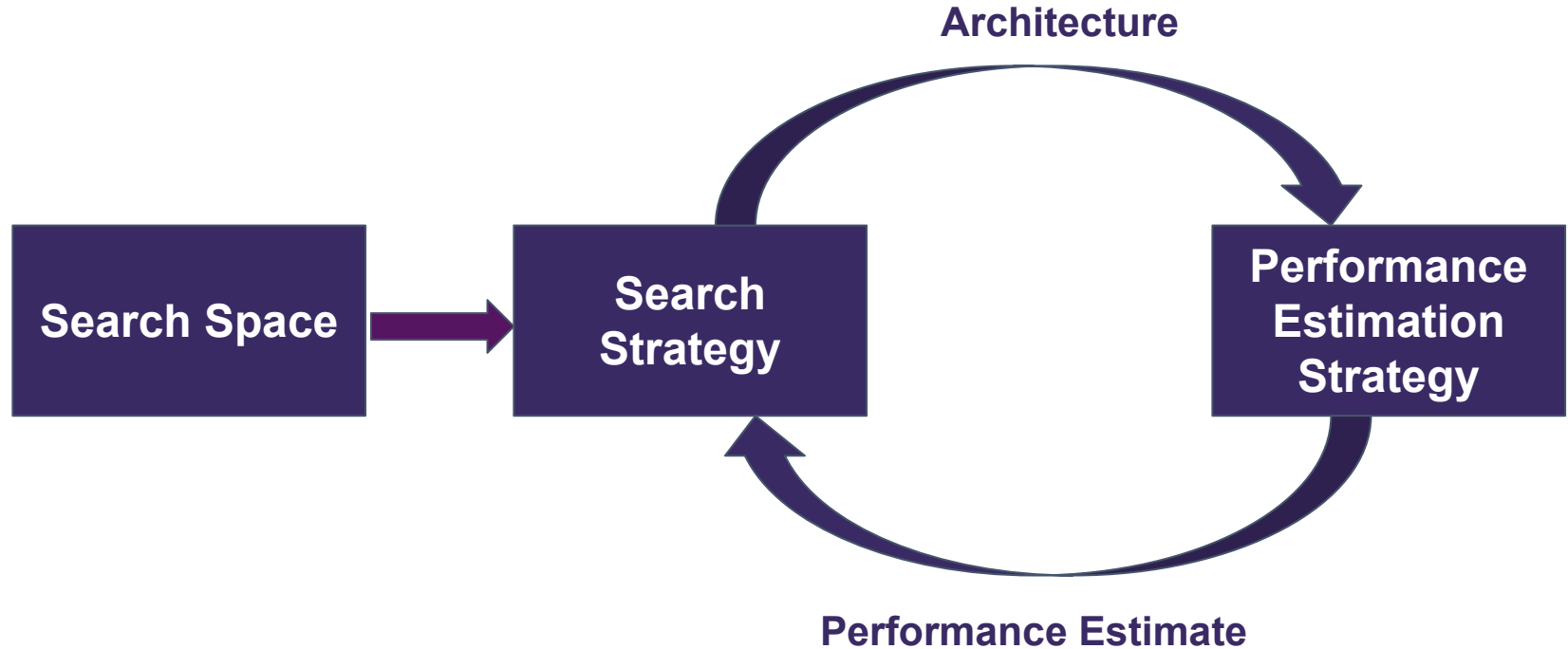
Neural Architecture Search (NAS)

NEURAL ARCHITECTURE SEARCH (NAS)

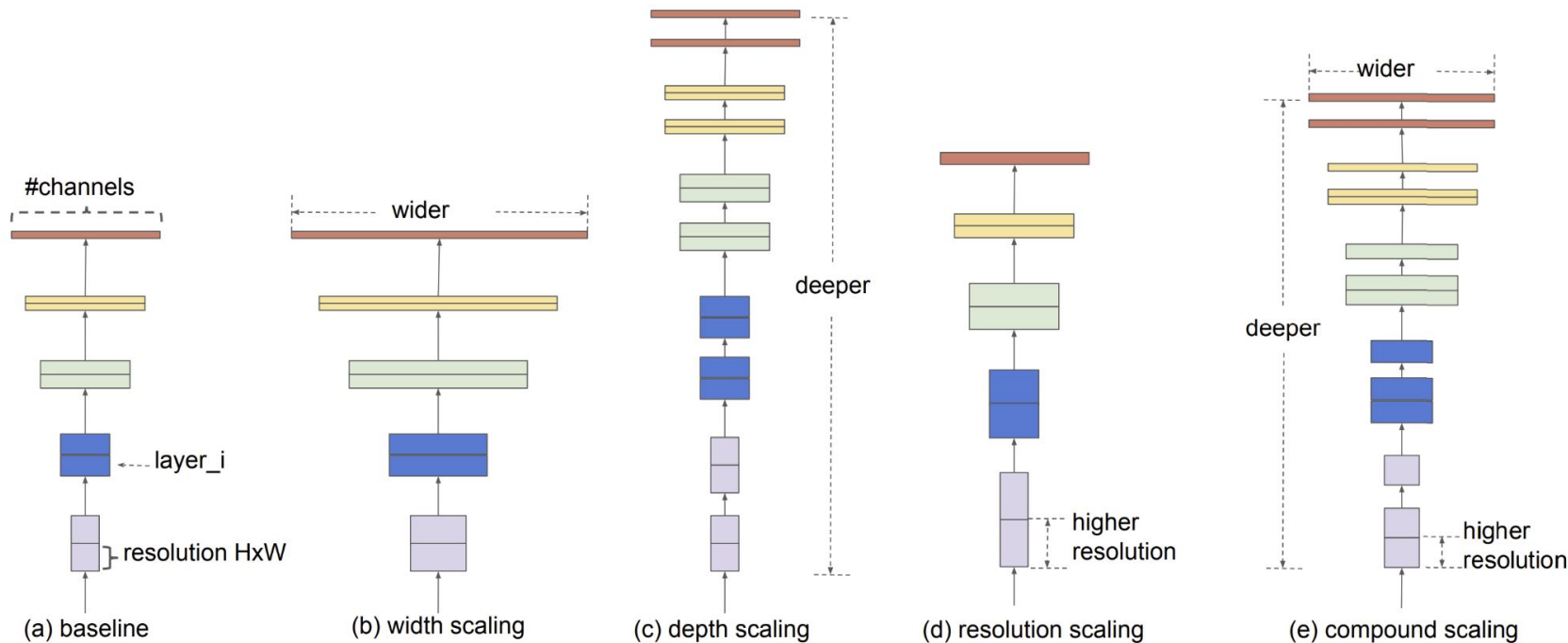
Search for the best neural network architecture on different hierarchical levels given a task.



The Components of NAS



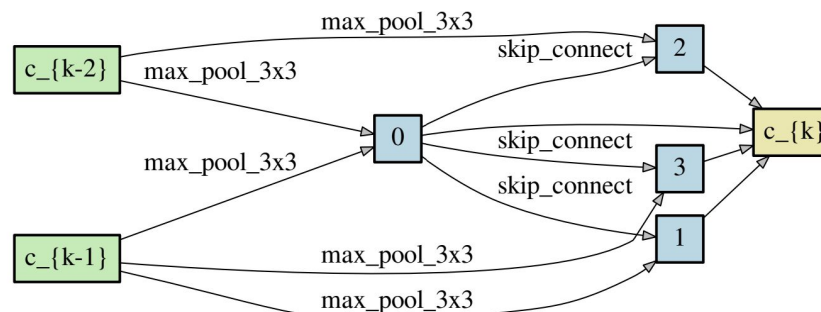
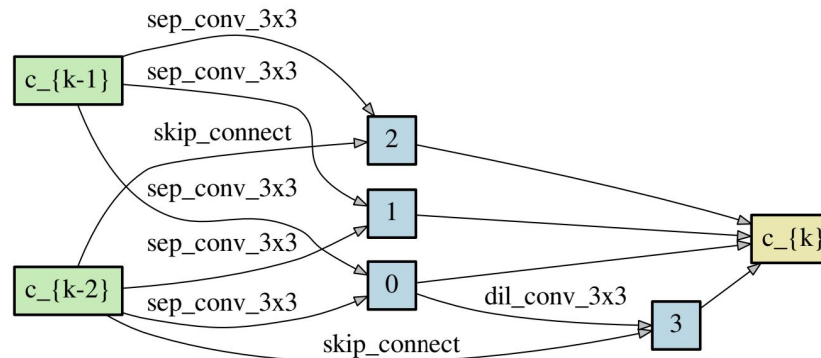
Search Space 1: Macro NAS



→ direct relationship to HPO: NAS as HPO

Source: [\[Tan & Le. 2019\]](#)

Search Space 2: Cell-based NAS



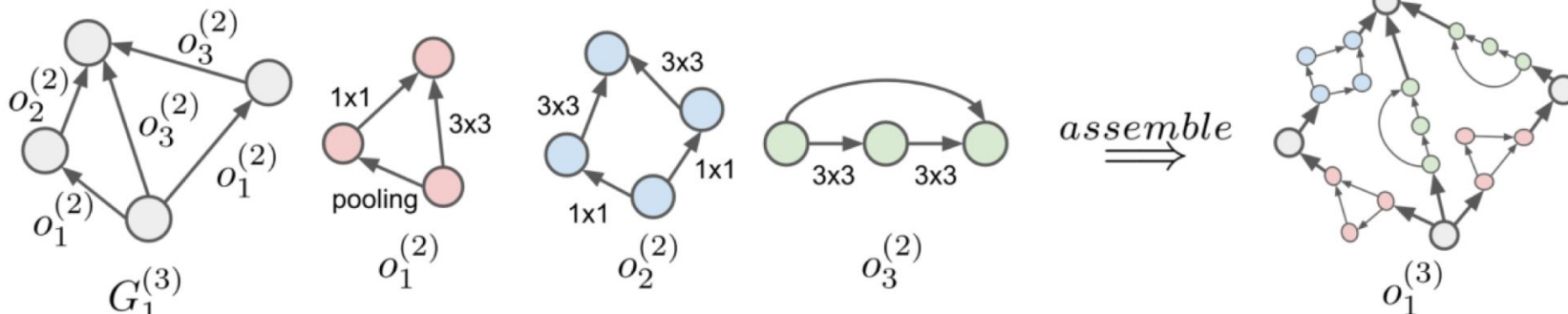
Source: [Liu et al. 2019]

Search Space 3: Hierarchical NAS

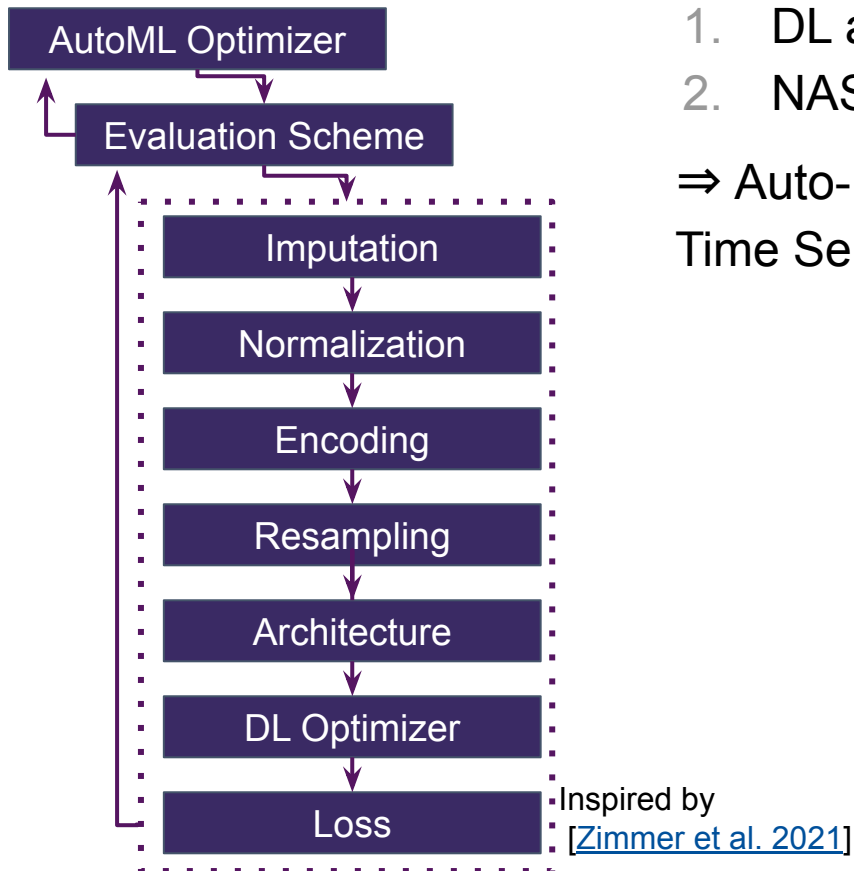
Search on multiple levels of the hierarchy

- **Lower levels:** create reusable building blocks
- **Higher levels:** combine building blocks

Like transformers are composed of lower-level building blocks (e.g., attention)



Source: [\[Liu et al, 2018\]](#)



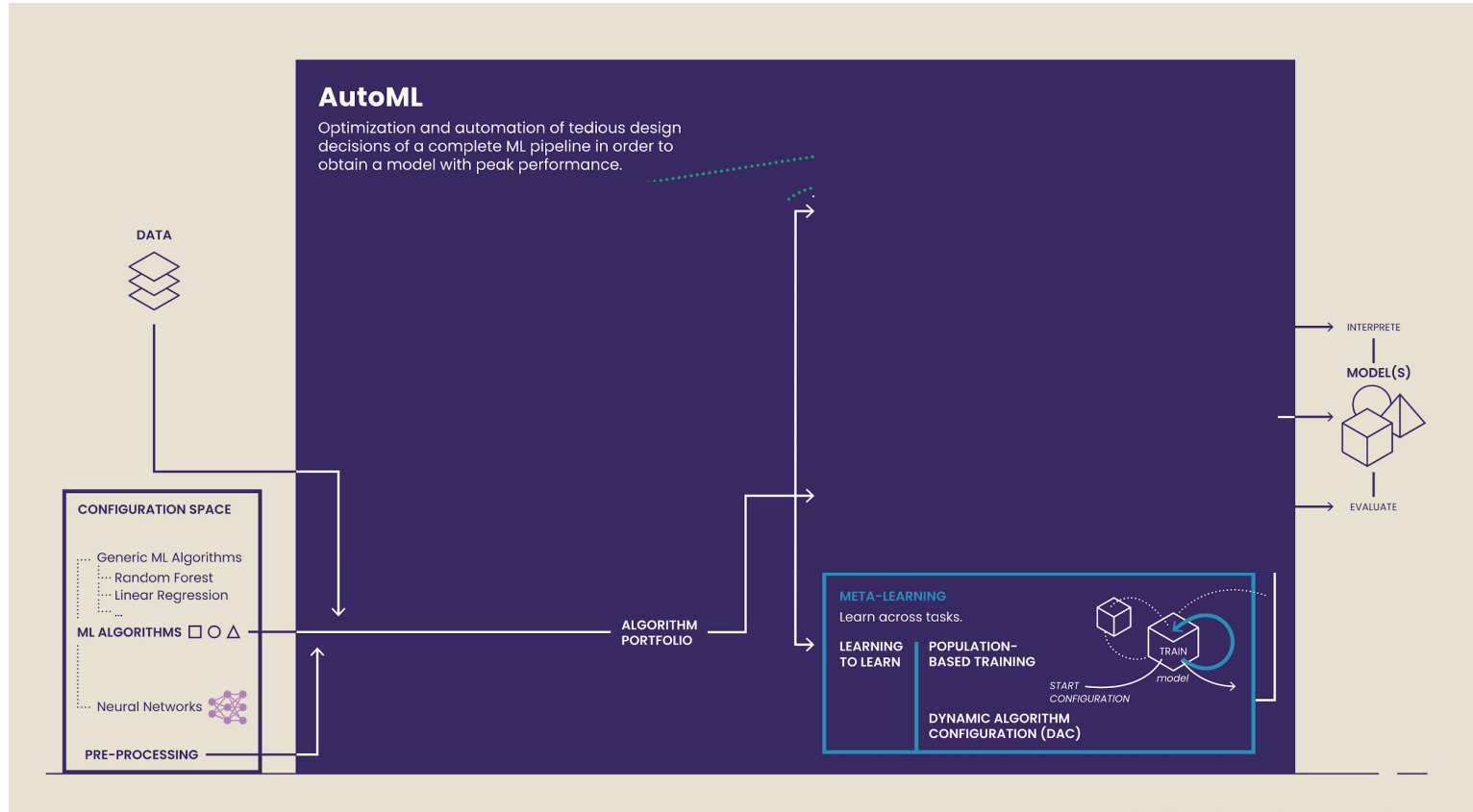
1. DL also includes complex pipelines
2. NAS & HPO need to go hand in hand

⇒ Auto-PyTorch [\[Zimmer et al. 2021\]](#) and Auto-PyTorch for Time Series Forecasting [\[Deng et al. 2022\]](#)

```
# initialise Auto-PyTorch api
api = TabularClassificationTask()

# Search for an ensemble of machine learning algorithms
api.search(
    X_train=X_train,
    y_train=y_train,
    X_test=X_test,
    y_test=y_test,
    optimize_metric='accuracy',
    total_walltime_limit=300,
    func_eval_time_limit_secs=50
)

# Calculate test accuracy
y_pred = api.predict(X_test)
```



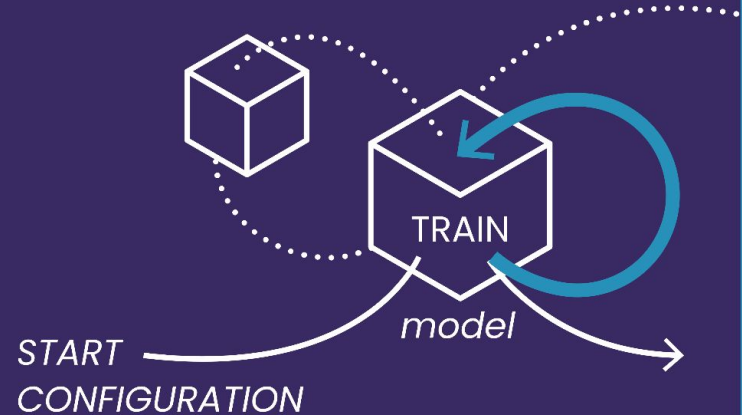
META-LEARNING

Learn across tasks.

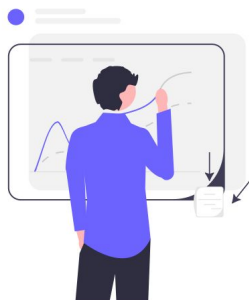
LEARNING
TO LEARN

POPULATION-
BASED TRAINING

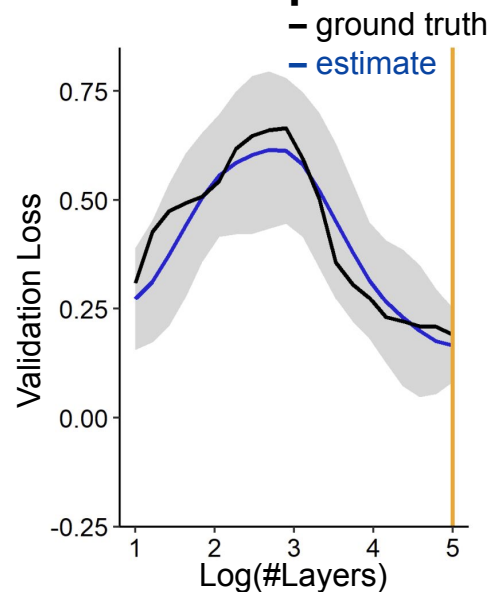
DYNAMIC ALGORITHM
CONFIGURATION (DAC)



Performance prediction

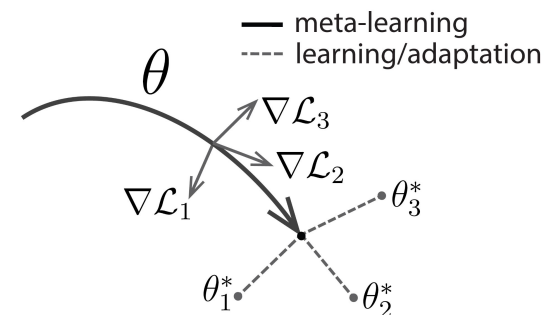


Hyperparameter Effects & Importance

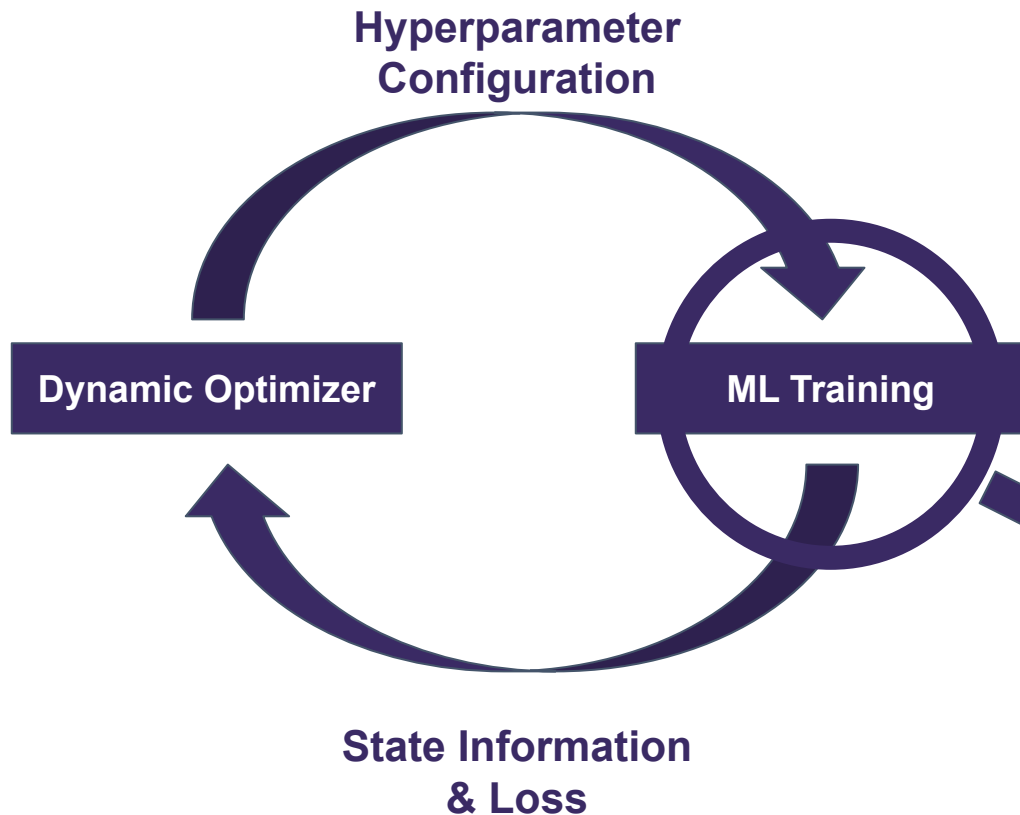


Source: [\[Moosbauer et al. 2021\]](#)

Learning NN weight initializations

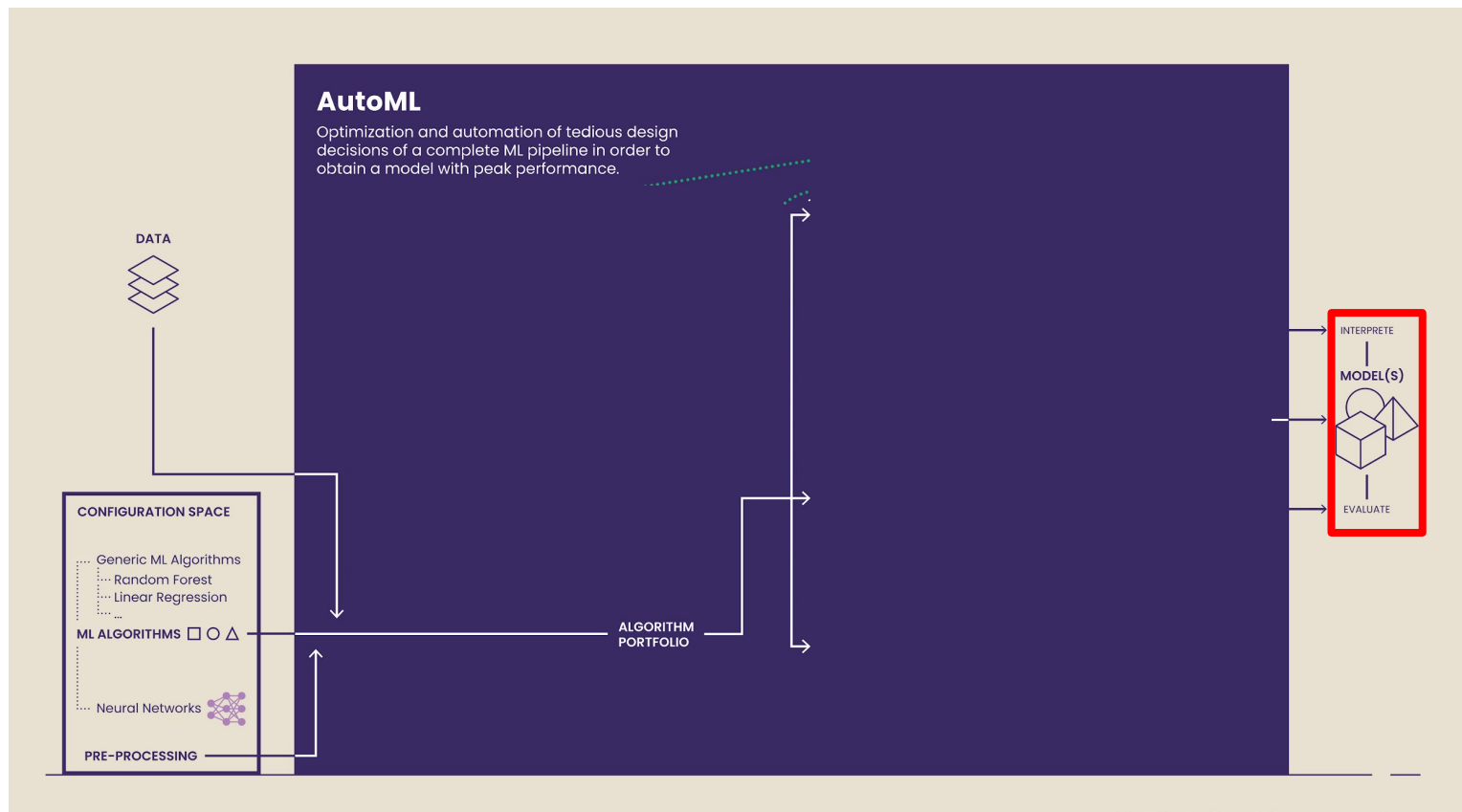


Source: [\[Finn et al. 2017\]](#)

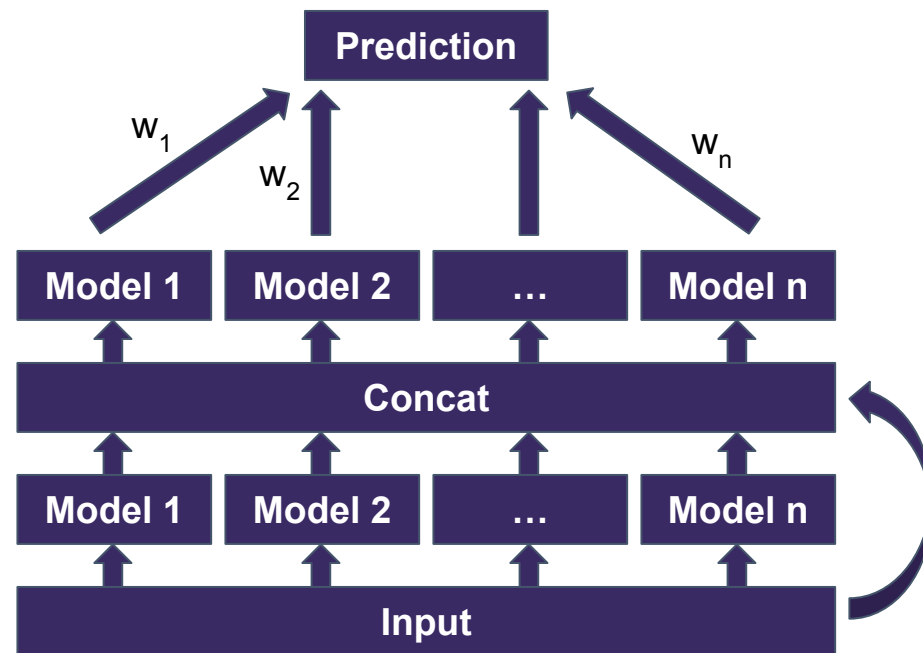
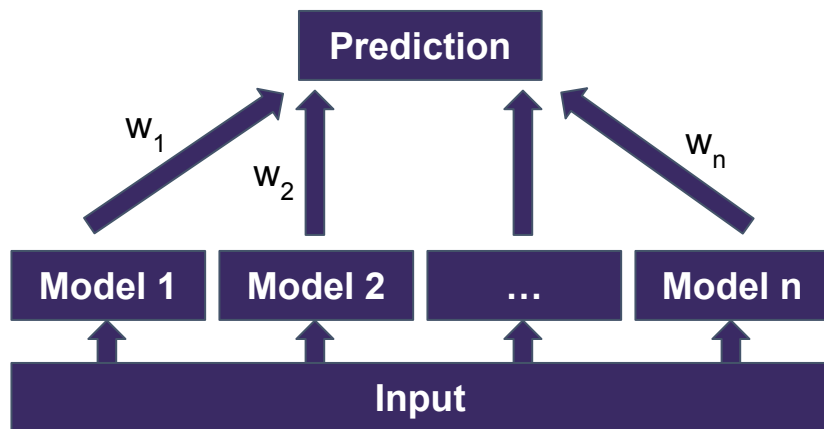


- Population-based Training
[[Jaderberg et al. 2017](#)]
- Population-based Bandits
[[Parker-Holder et al. 2020](#)]
- Dynamic Algorithm Configuration via RL
[[Biedenkapp et al. 2020](#),
[Adriaensen et al. 2022](#)]
- Adapting Bayesian Optimization
[Benjamins et al. 2022]

Final Step of AutoML

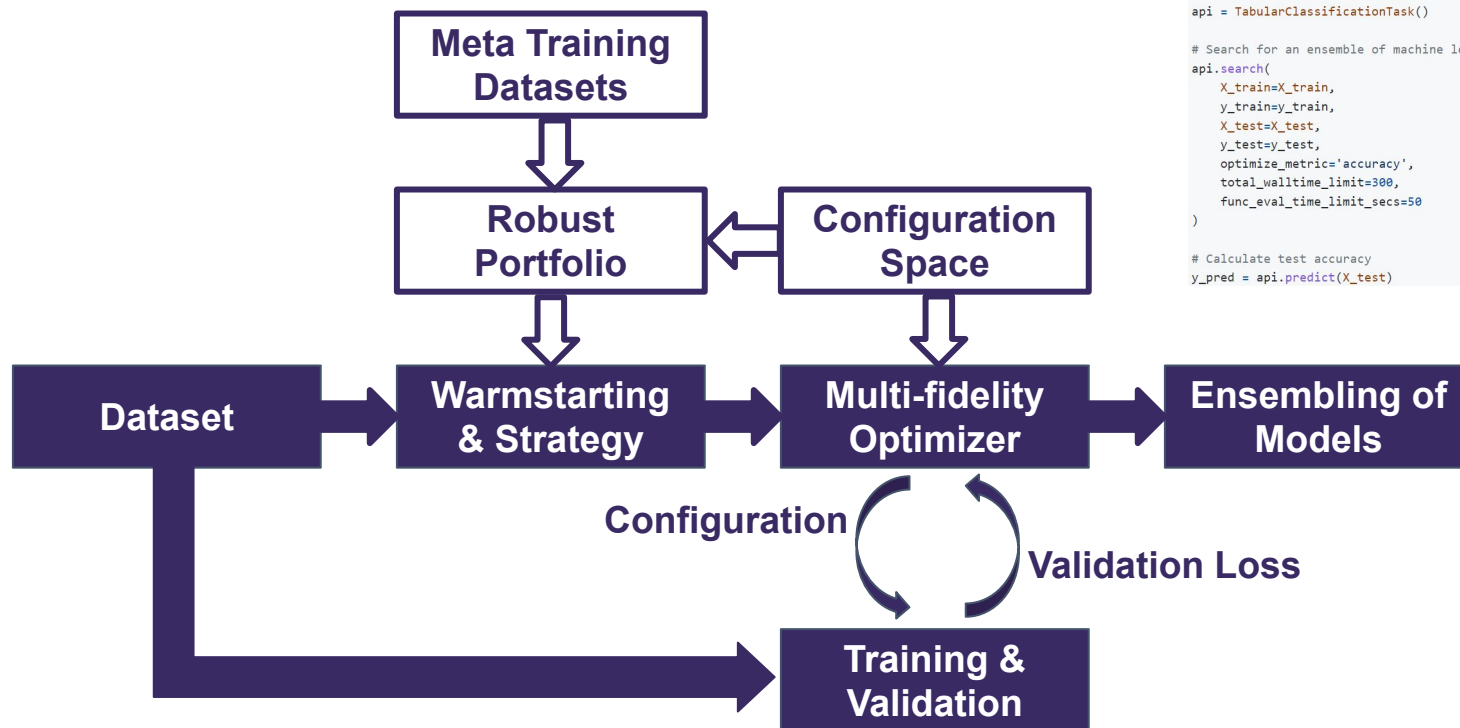


Ensembling vs Stacking



Source [Erickson et al. 2020]

Auto-Sklearn [Feurer et al. 2015, Feuer et al. 2022] & Auto-PyTorch [Zimmer et al. 2021]

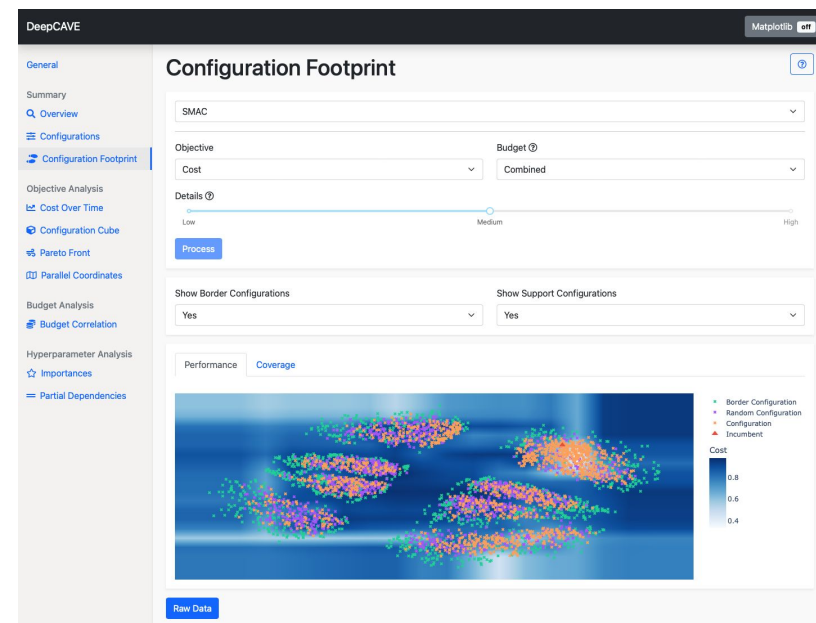
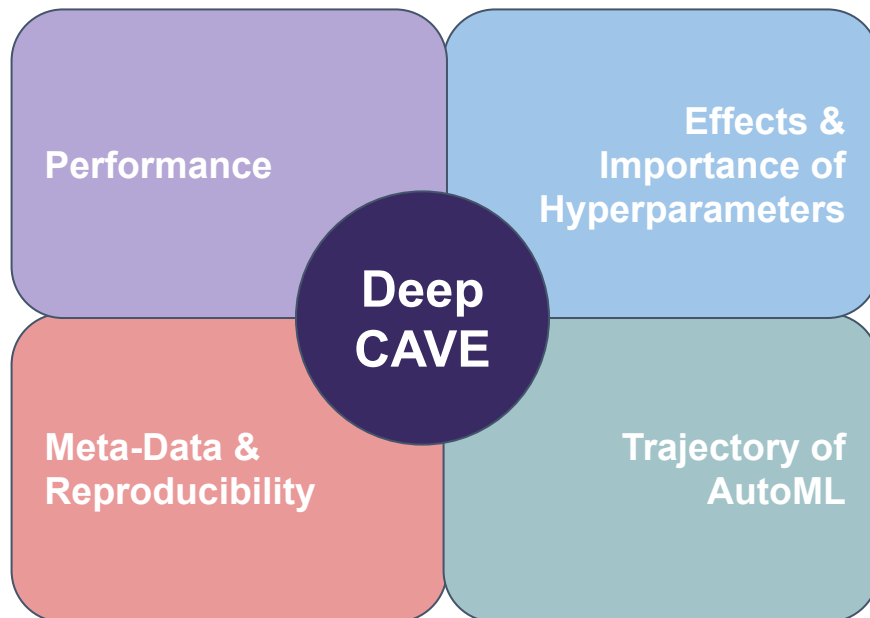


```
# initialise Auto-PyTorch api
api = TabularClassificationTask()

# Search for an ensemble of machine learning algorithms
api.search(
    X_train=X_train,
    y_train=y_train,
    X_test=X_test,
    y_test=y_test,
    optimize_metric='accuracy',
    total_walltime_limit=300,
    func_eval_time_limit_secs=50
)

# Calculate test accuracy
y_pred = api.predict(X_test)
```

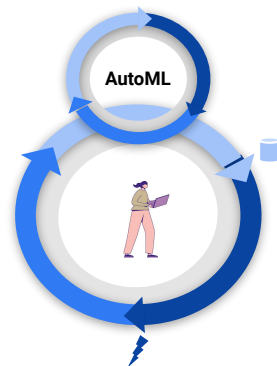
Monitoring AutoML [Sass et al. 2022]



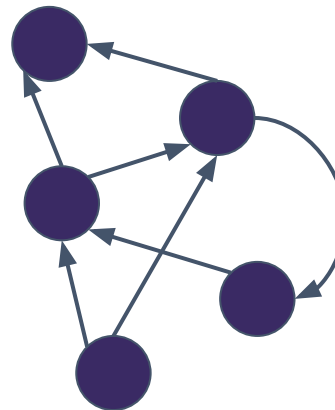
Selection of Open Challenges



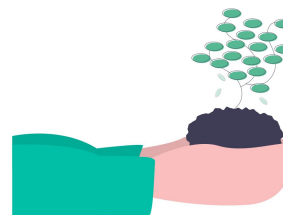
**Scaling up AutoML
for very large models**



Human-centered AutoML

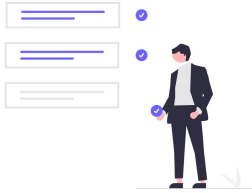


**Finding substantially
novel architectures**

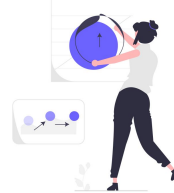


Green AutoML

Are Data Scientists still needed? Yes



**Determine your
objectives, metrics
and constraints**



**Design the
configuration space**



**Bring in the domain
knowledge**



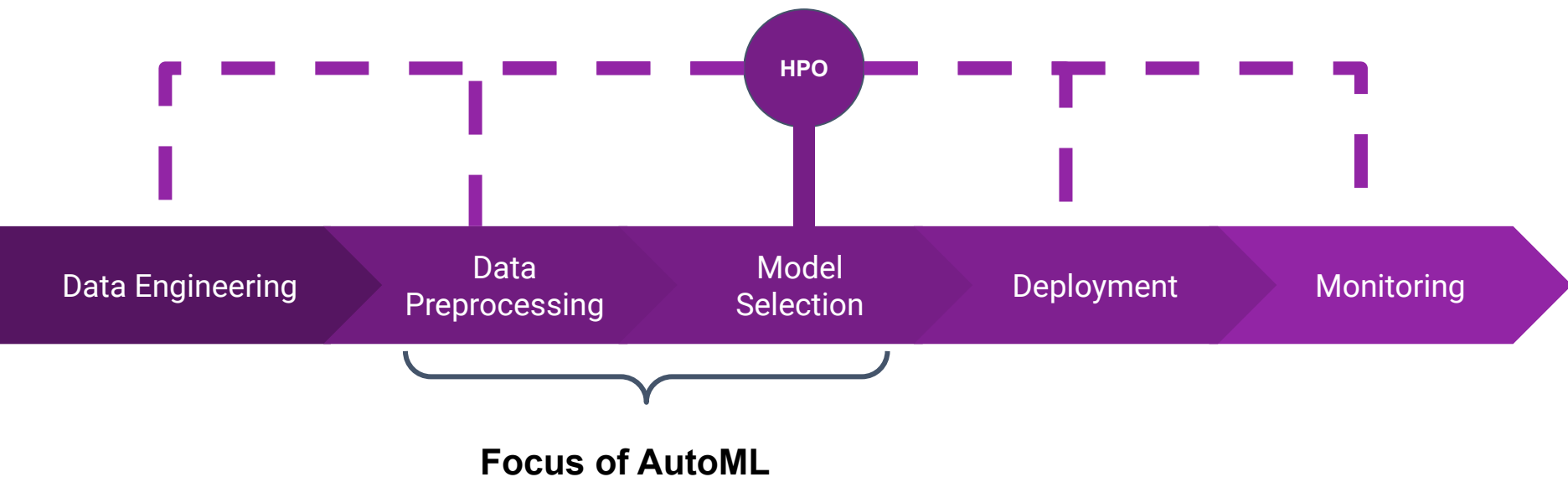
Determine Budgets

Running AutoML



Monitor AutoML

HPO → AutoML → AutoDS



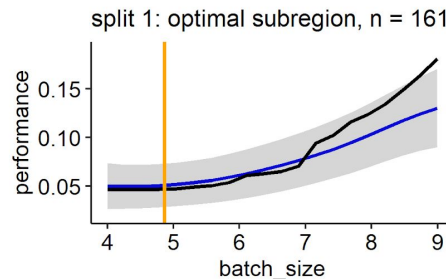
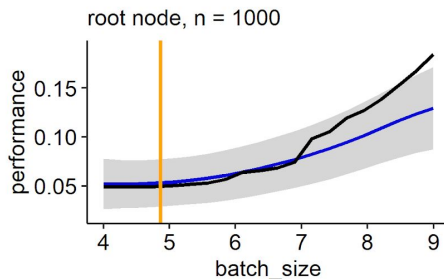
Can we explain what AutoML figured out?

[[Moosbauer et al. NeurIPS'21](#), [Moosbauer et al. 2022](#)]

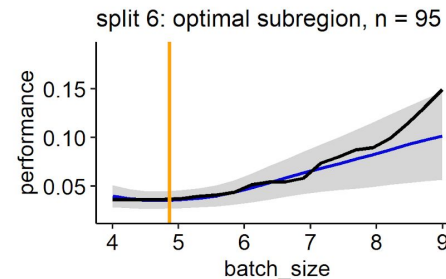


Explaining Hyperparameter Effects via PDPs

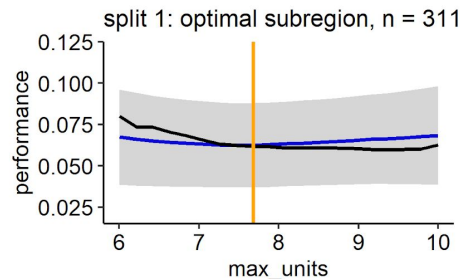
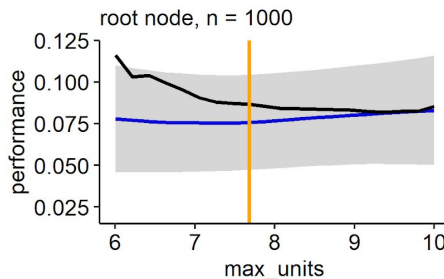
Ground truth
PDP
incumbent



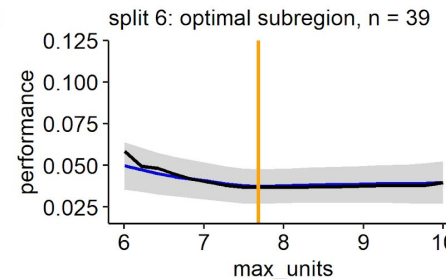
Subregion definition:
weight_decay <= 0.086



Subregion definition:
num_layers <= 4.5,
weight_decay <= 0.0178,
max_dropout <= 0.6966



Subregion definition:
batch_size <= 7.5329



Subregion definition:
max_dropout <= 0.7305,
num_layers <= 4.5,
batch_size <= 6.1739,
weight_decay <= 0.0172

For, a subset S of the hyperparameters, the partial dependence function is:

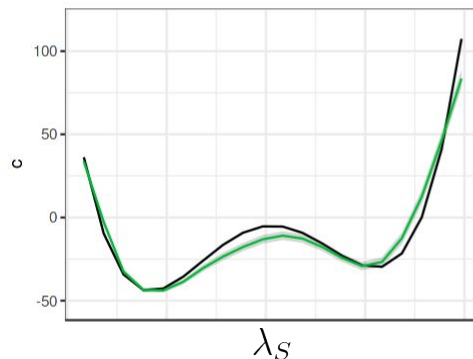
$$c_S(\lambda_S) := \mathbb{E}_{\lambda_C} [c(\lambda)] = \int_{\Lambda_C} c(\lambda_S, \lambda_C) d\mathbb{P}(\lambda_C)$$

and can be approximated by Monte-Carlo integration on a surrogate model:

$$\hat{c}_S(\lambda_S) = \frac{1}{n} \sum_{i=1}^n \hat{m}(\lambda_S, \lambda_C^{(i)})$$

where $\left(\lambda_C^{(i)}\right)_{i=1, \dots, n} \sim \mathbb{P}(\lambda_C)$ and λ_S for a set of grid points.

→ Average of ICE curves.



Green: PDP

Black: Ground truth

[[Hutter et al. 2014](#)] showed how to do this efficiently for RFs as surrogate models.

Partial Dependence Plots with Uncertainties

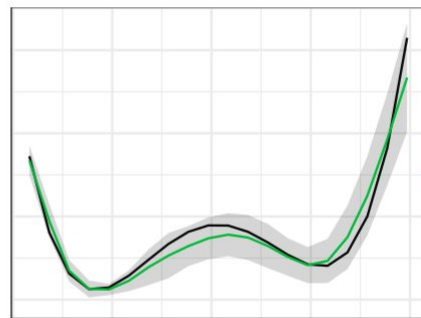
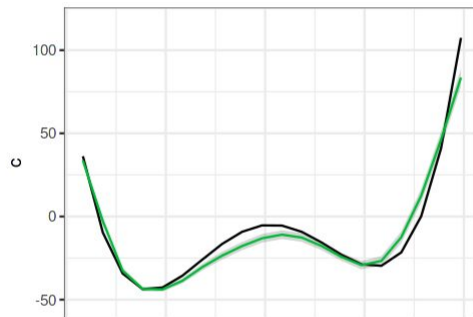
$$\begin{aligned} & \hat{s}_S^2(\lambda_S) \\ &= \mathbb{V}_{\hat{c}} [\hat{c}_S(\lambda_S)] \\ &= \mathbb{V}_{\hat{c}} \left[\frac{1}{n} \sum_{i=1}^n \hat{c}(\lambda_S, \lambda_C^{(i)}) \right] \\ &= \frac{1}{n^2} \mathbf{1}^\top \hat{K}(\lambda_S) \mathbf{1}. \end{aligned}$$

→ requires a kernel correctly specifying the covariance structure (e.g., GPs).

Approximation:

$$\hat{s}_S^2(\lambda_S) \approx \frac{1}{n} \sum_{i=1}^n \hat{K}(\lambda_S)_{i,i}$$

→ Model-agnostic (local) approximation

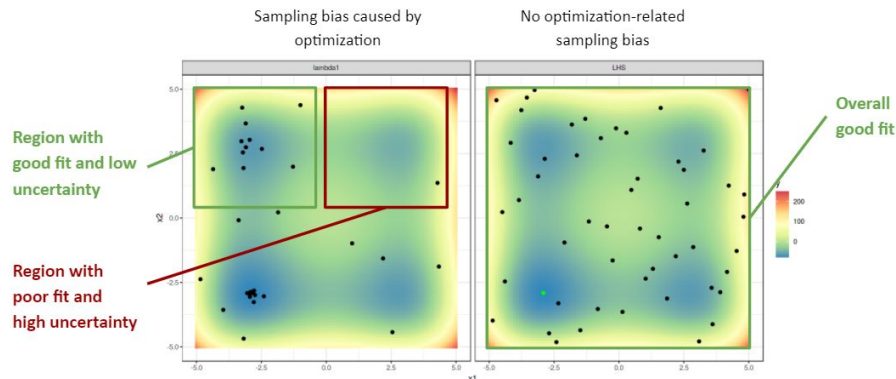
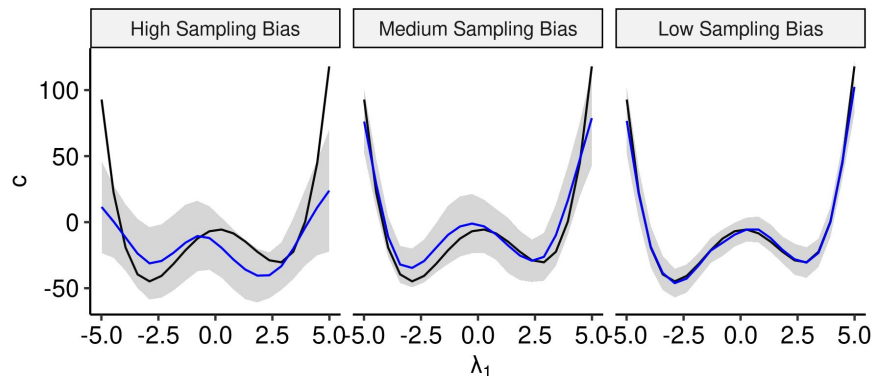


Ground truth
PDP
Uncertainty

Impact of Sampling Bias in Explaining AutoML

- Simply using all observations from AutoML tools might lead to misleading PDPs
- Uncertainty estimates help to quantify the poor fits

→ of course, sampling bias is wanted and the solution cannot be to change the sampling behavior

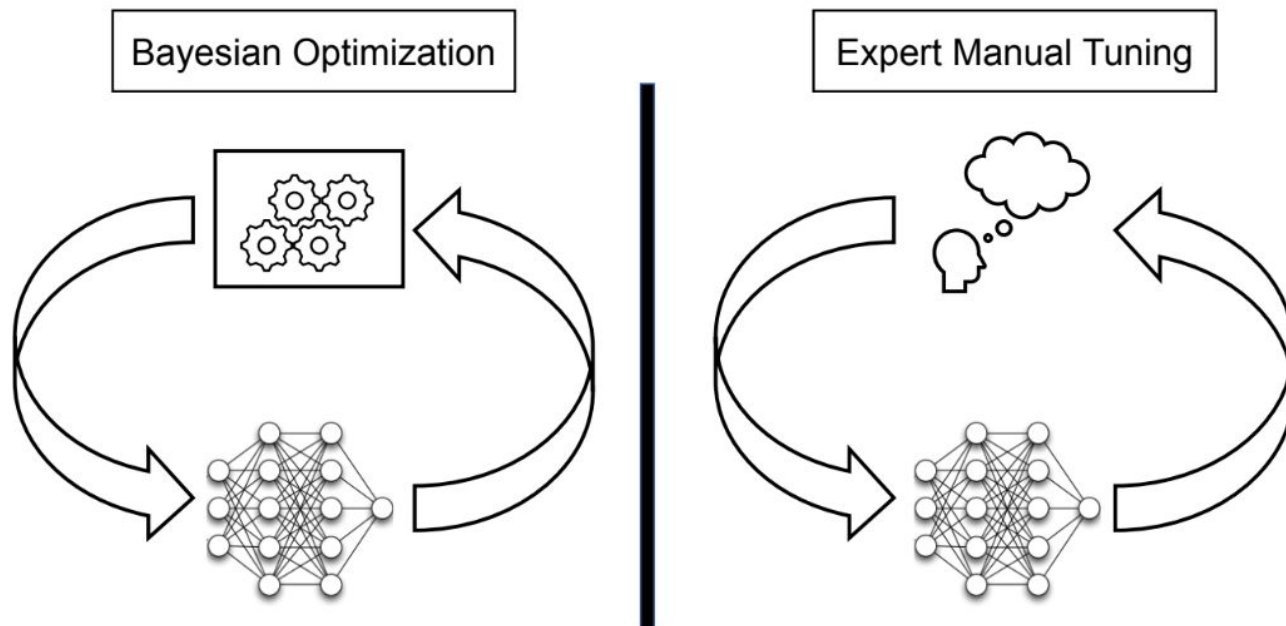


Can AutoML consider expert knowledge?

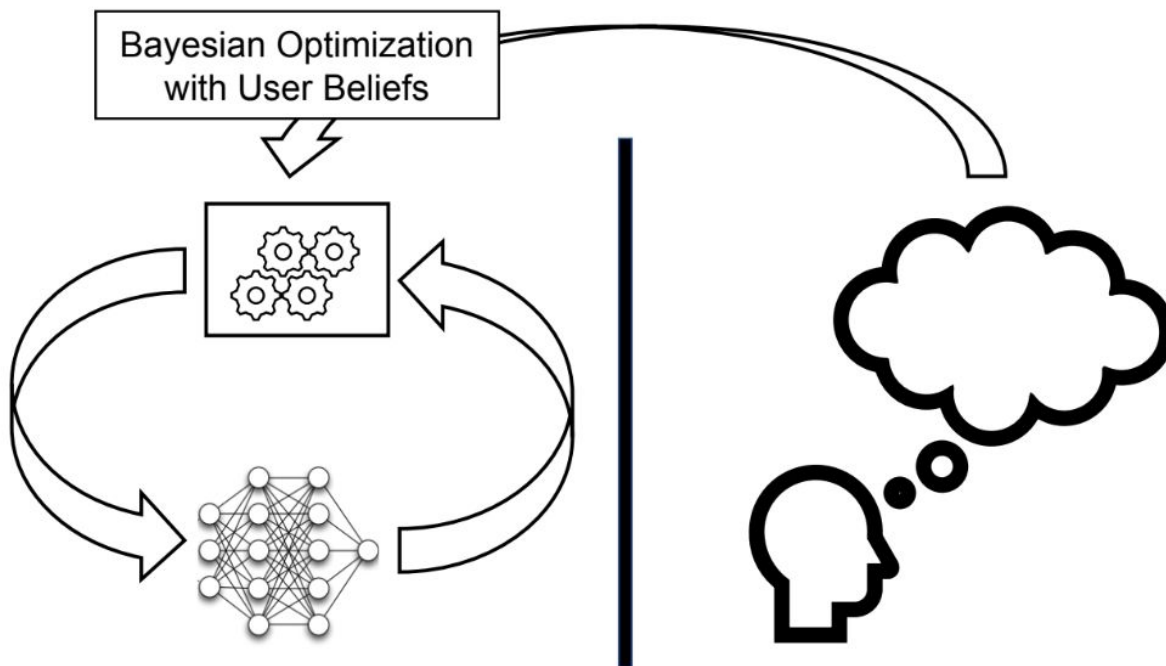
[Hvarfner et al. ICLR'22]

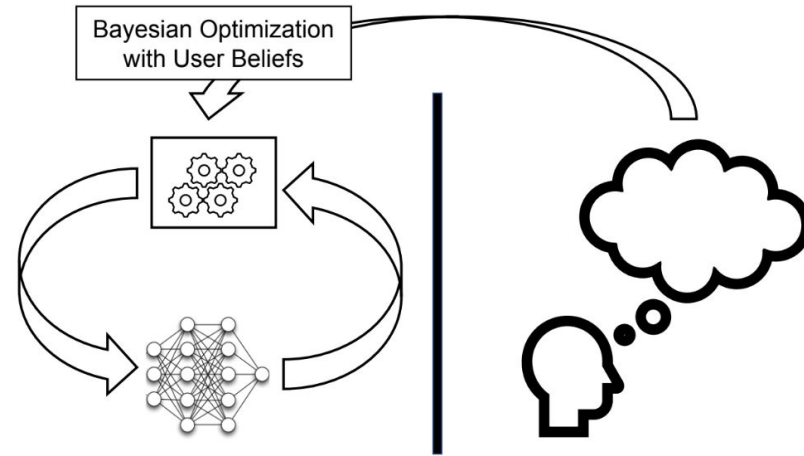


Bayesian Optimization vs Manual Tuning



Bayesian Optimization with Expert Knowledge



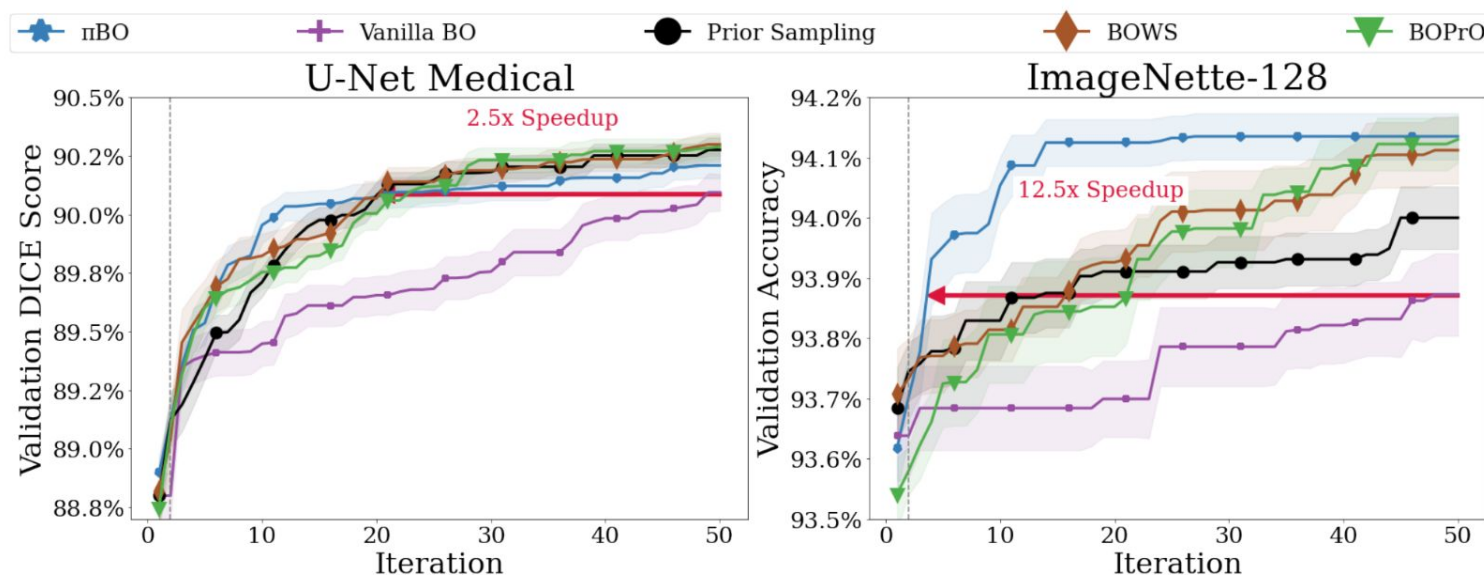


$$\mathbf{x}_n \in \arg \max_{\mathbf{x} \in \mathcal{X}} \alpha(\mathbf{x}, \mathcal{D}_n) \pi(\mathbf{x})^{\beta/n}$$

Acquisition Function

User Prior

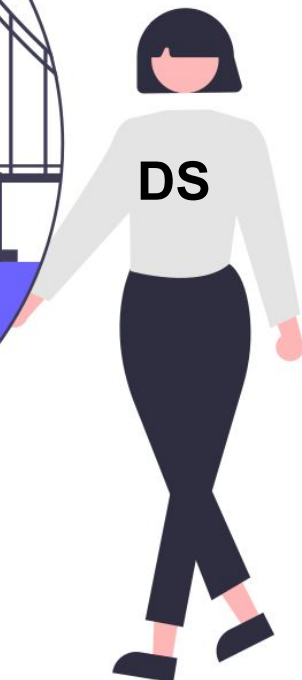
Speed of forgetting user prior



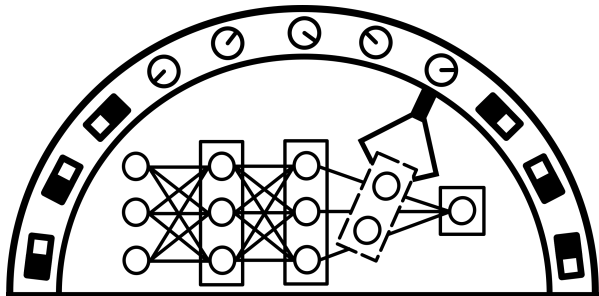
- Uses expert knowledge to speed up Bayesian Optimization
- Robust also against wrong believes
- Substantially speeds up AutoML

Will AutoML replace Data Scientists?

Application



**AutoML: Helping to
bridge application
and data science.**



AutoML.org



/AutoML_org/



/automl/



<https://tinyurl.com/automlyt>

Funded by:



European Research Council
Established by the European Commission



Deutsche
Forschungsgemeinschaft



Federal Ministry
of Education
and Research



Federal Ministry
for Economic Affairs
and Energy



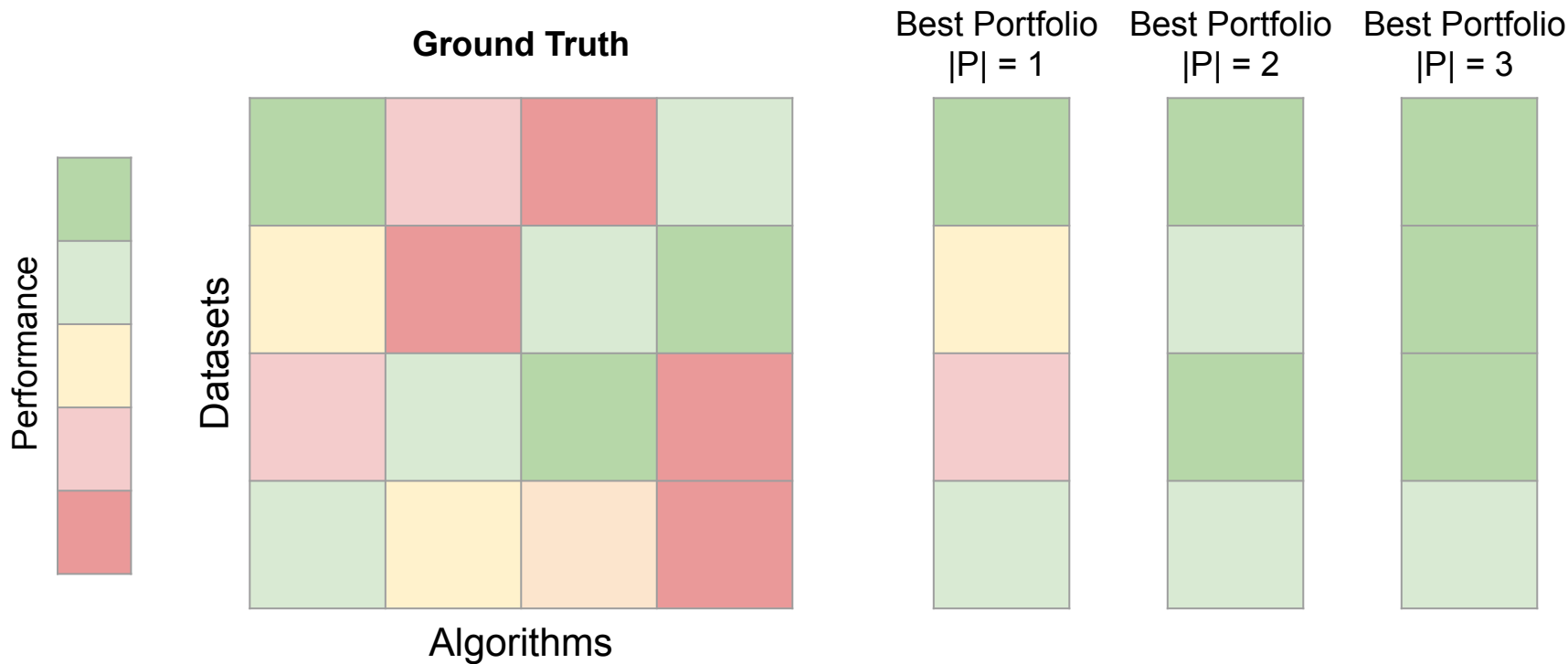
1 1
1 0 2
1 0 0 4

Leibniz
Universität
Hannover

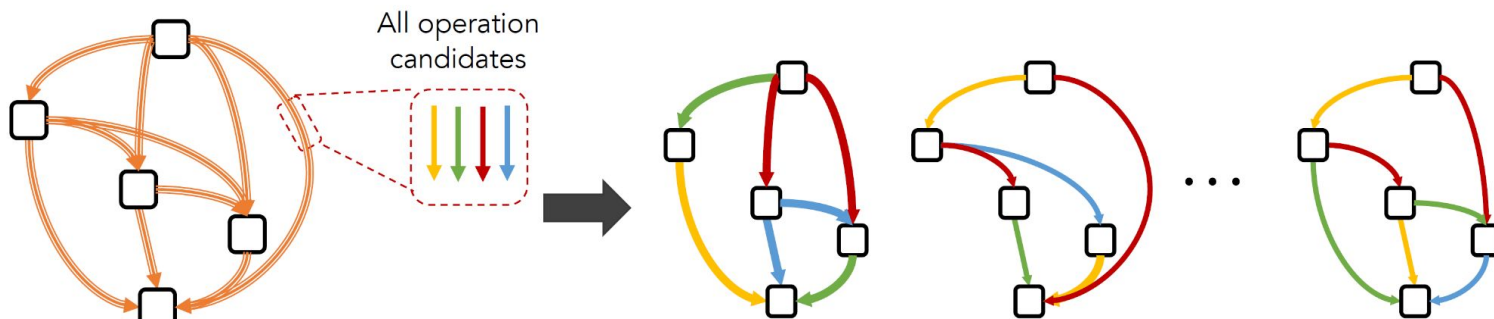


Backup slides

Portfolios for Warmstarting [\[Feurer et al. 2022\]](#)



Oneshot NAS: Weight Sharing Across Architectures



- For each choice between operations, the supernet includes all of them
- A linear number of weights shared by an exponential number of architectures
- Thus, updating the weights of one architecture simultaneously updates parts of the weights of exponentially many other architectures



Very hot topic in NAS, but no consistent improvements over trivial baselines, such as #parameters or FLOPs

ZC proxies are a particular type of performance predictor

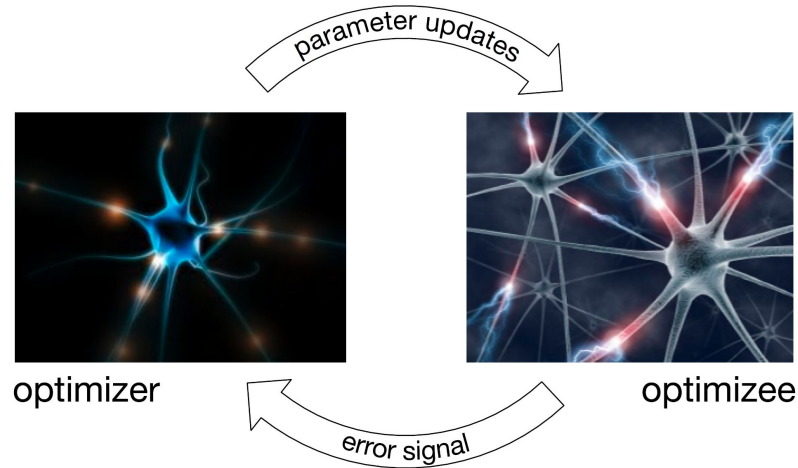
- They aim to judge the performance of an architecture in a few seconds
- Often by a single forward pass on a mini-batch
- Thus, the term “zero-cost”

Examples

- Change of error when dropping network weights
- Dissimilarity of activation patterns for points in a batch

Very hot topic in NAS, but no consistent improvements over using number of parameters or FLOPS

E.g., “Learning to learn by gradient descent by gradient descent” [[Chen et al. 2016](#)]



Source: [[Chen et al. 2016](#)]

E.g., Alpha-Zero [[Silver et al. 2017](#)]

Maturity of AutoML

