



Deep Learning Optimisé - Jean Zay

Optimisation des hyperparamètres



INSTITUT DU
DÉVELOPPEMENT ET DES
RESSOURCES EN
INFORMATIQUE
SCIENTIFIQUE



HPO = Hyperparameter Optimisation

Hyperparameters ◀

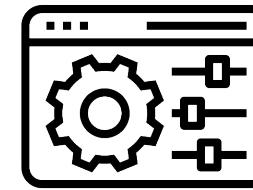
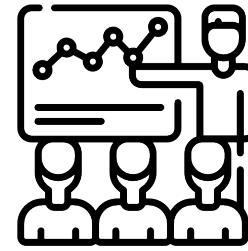
HPO ◀

Related Problems ◀

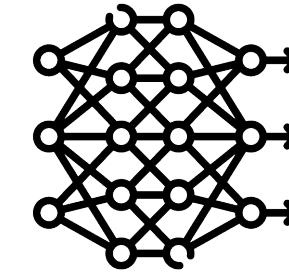
Hyperparameters

In machine learning, a hyperparameter is **a parameter whose value is used to control the learning process**.
By contrast, the values of other parameters (typically node weights) are derived via training.

Hyperparameters



Parameters

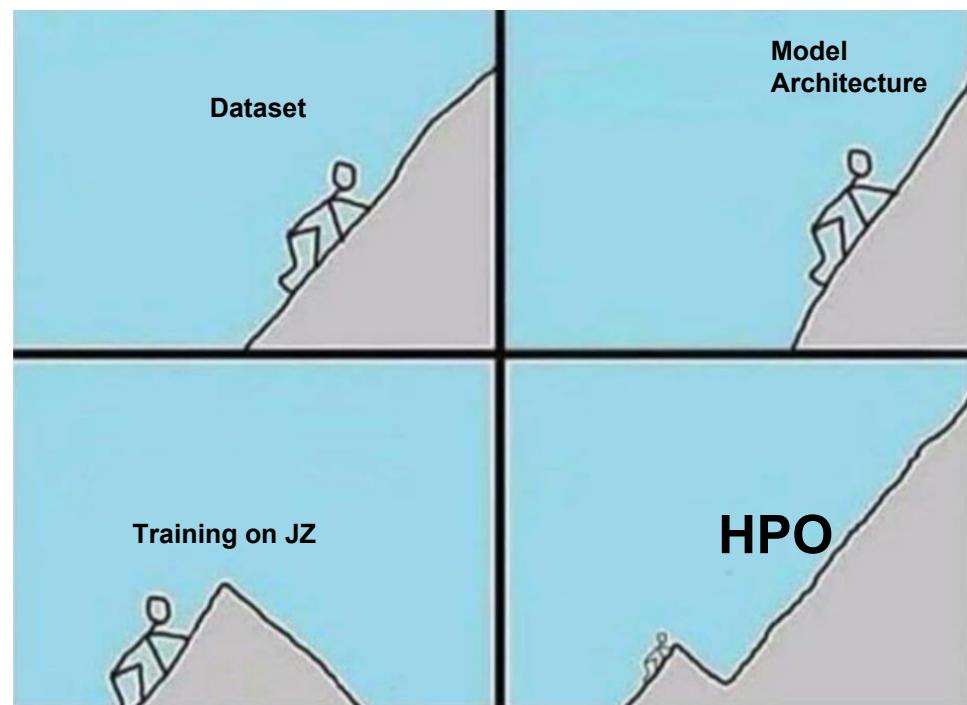
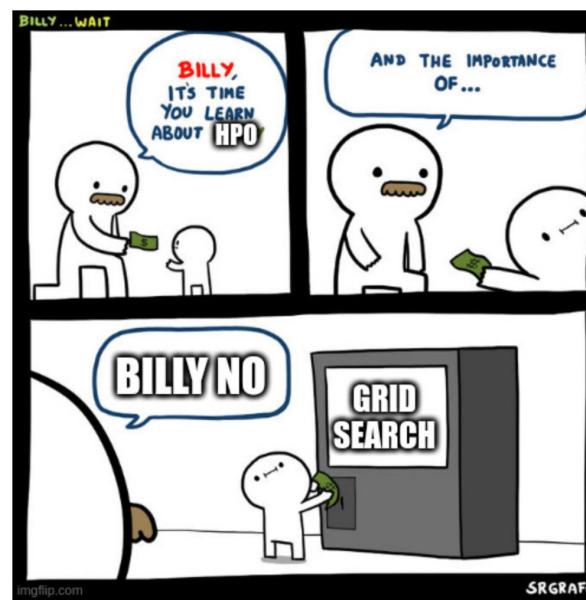


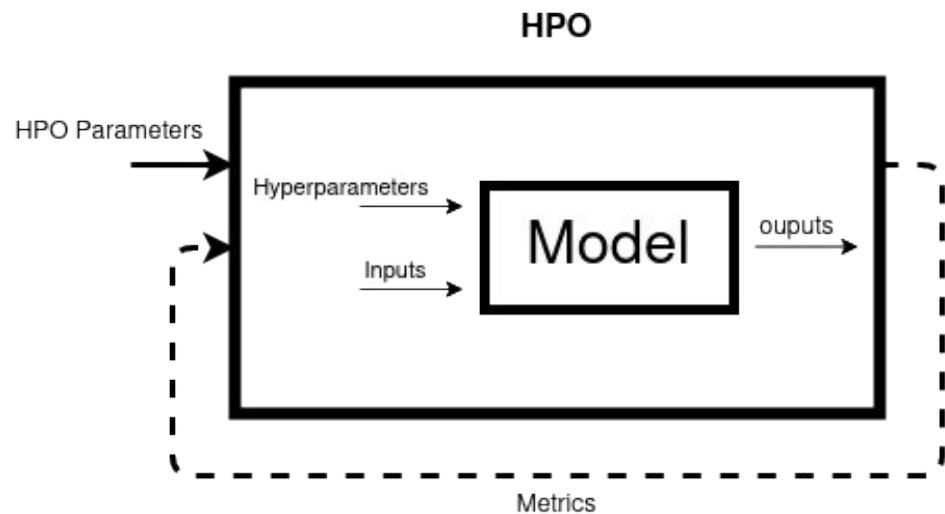
HPO : Hyperparameter Optimisation

Machine learning algorithms are highly configurable by their hyperparameters.

These parameters often substantially influence the complexity, behavior, speed as well as other aspects of the learner, and their values must be selected with care in order to achieve optimal performance.

Human trial-and-error to select these values is time-consuming, often somewhat biased, error-prone and computationally irreproducible.



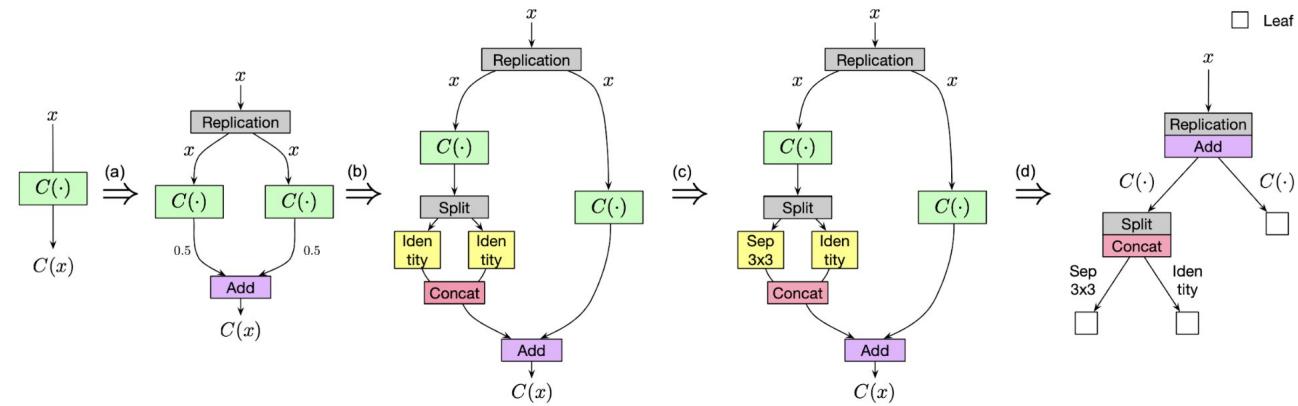
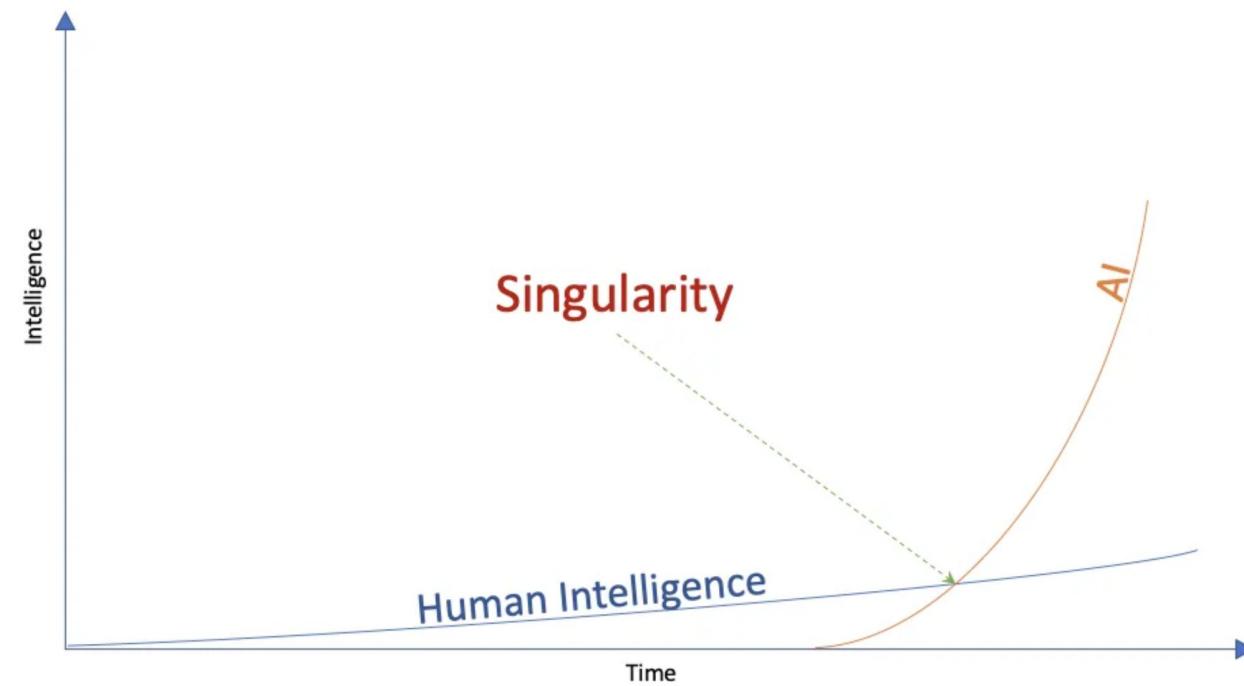


Hyperparameter Optimization == Bi-Level optimization problem



Related Problems

- Neural Architecture Search (**NAS**)
- Algorithm Selection and traditional Meta-Learning
- Algorithm configuration (**AC**)
- Dynamic Algorithm Configuration (**DAC**)
- Learning to learn and to optimize



A Comprehensive Survey of Neural Architecture Search: Challenges and Solution (<https://arxiv.org/pdf/2006.02903.pdf>)



Fastest wheel change on a moving car - Guinness World Records

Search Algorithms / Samplers

Basic ◀

Manual, Grid Search, Random Search

Bayesian Optimisation ◀

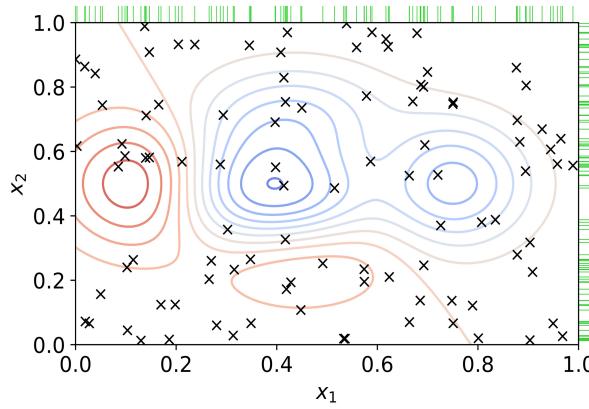
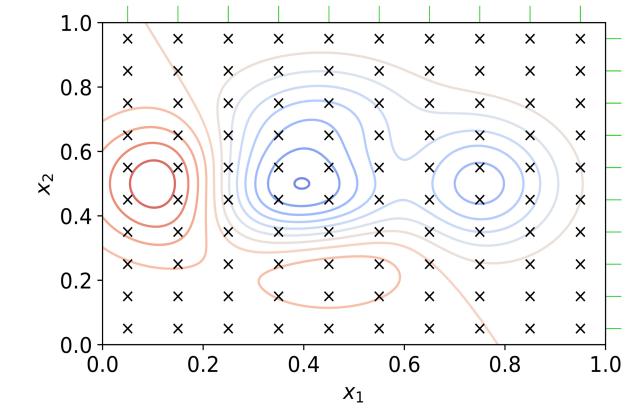
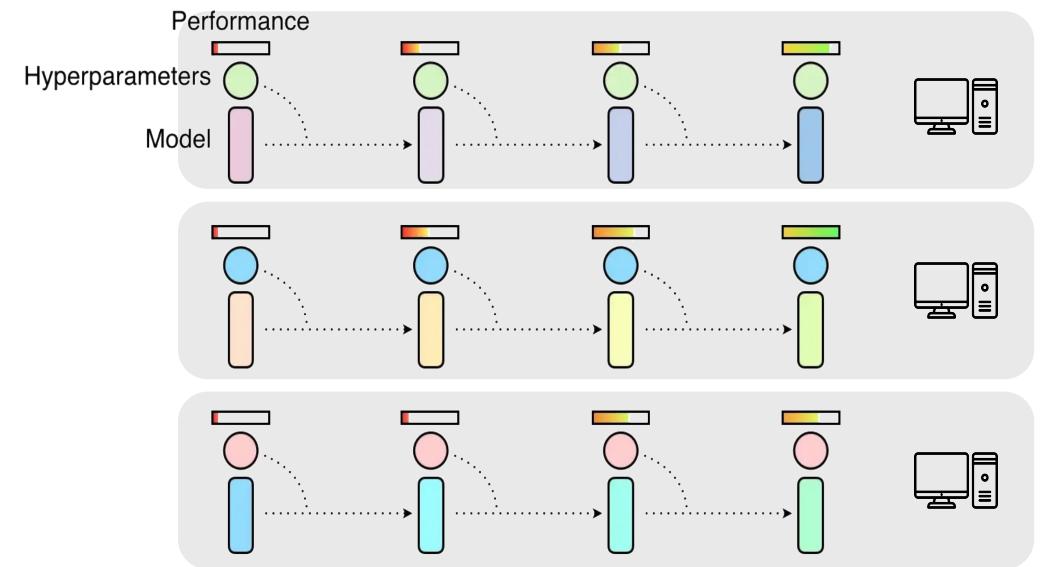
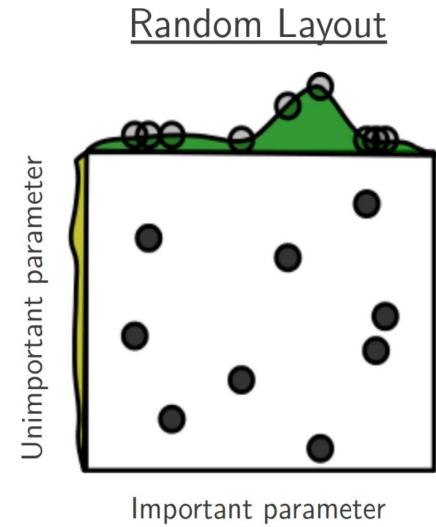
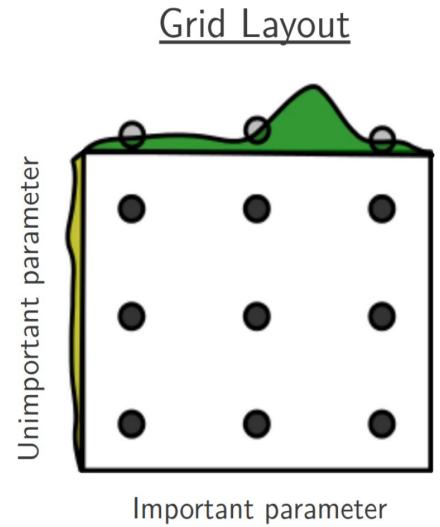
Tree-structured Parzen Estimator, Gaussian Process

Heuristic ◀

Genetic Algorithm, Particle Swarm Optimization

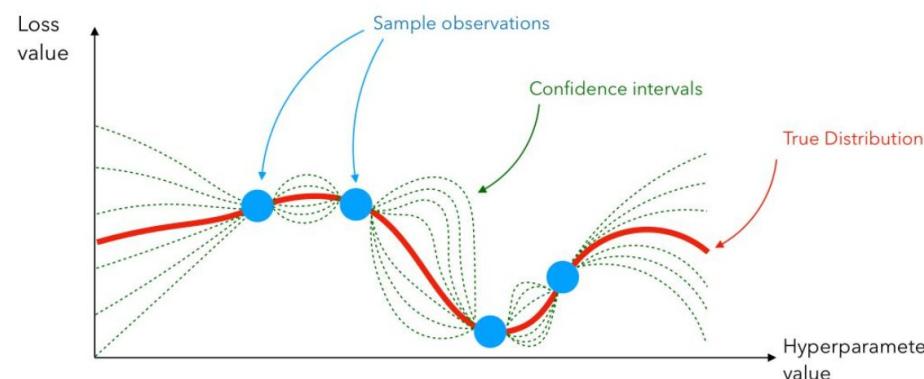
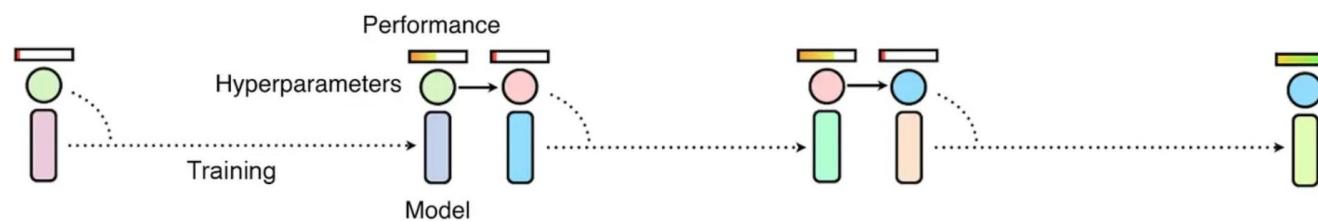
Gradient-based Optimization ◀

Basic : Grid & Random Search

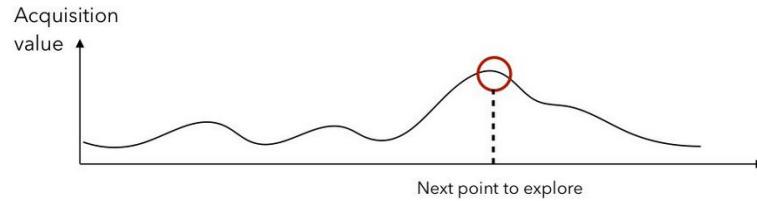


- Independent tests (which can be parallelized) which test a combination of hyperparameters.
- Very costly in resources and no guarantee of improved results.
- Random search is better for high dimensional space

Bayesian Optimization : TPE & GP

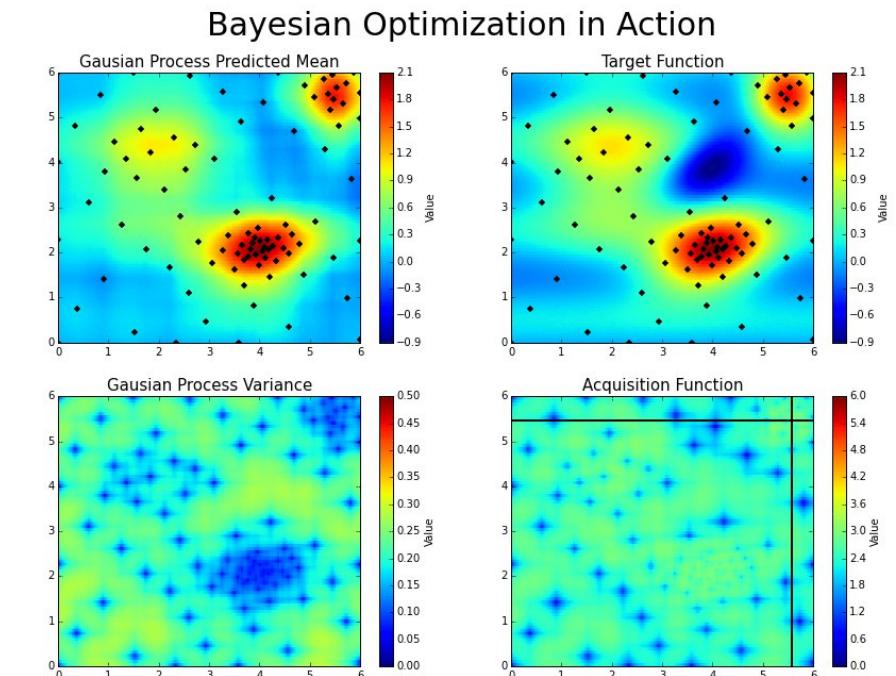


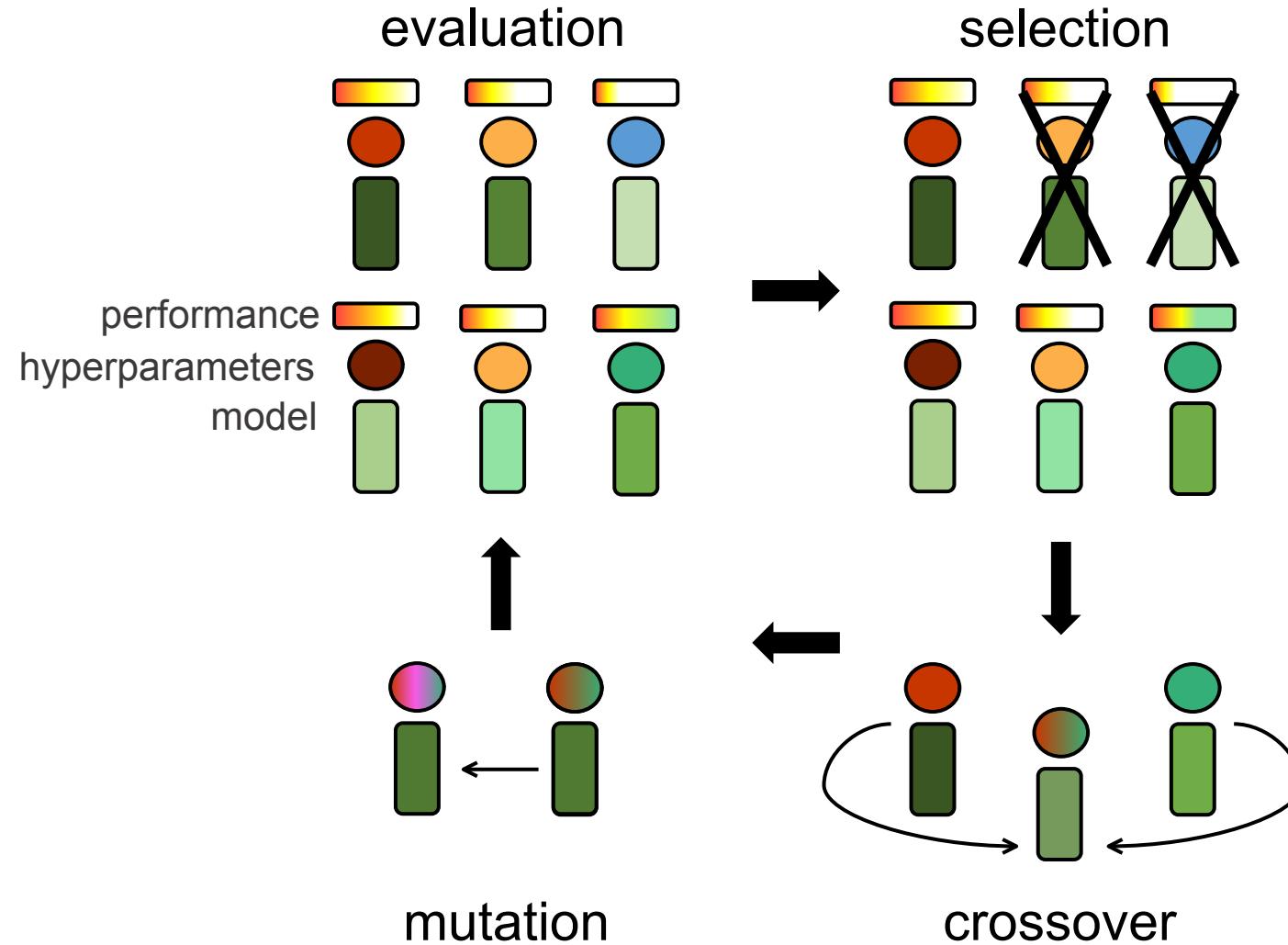
Expected metric score according to Hyper-parameters



Maximize Acquisition function e.g. Expected Improvement

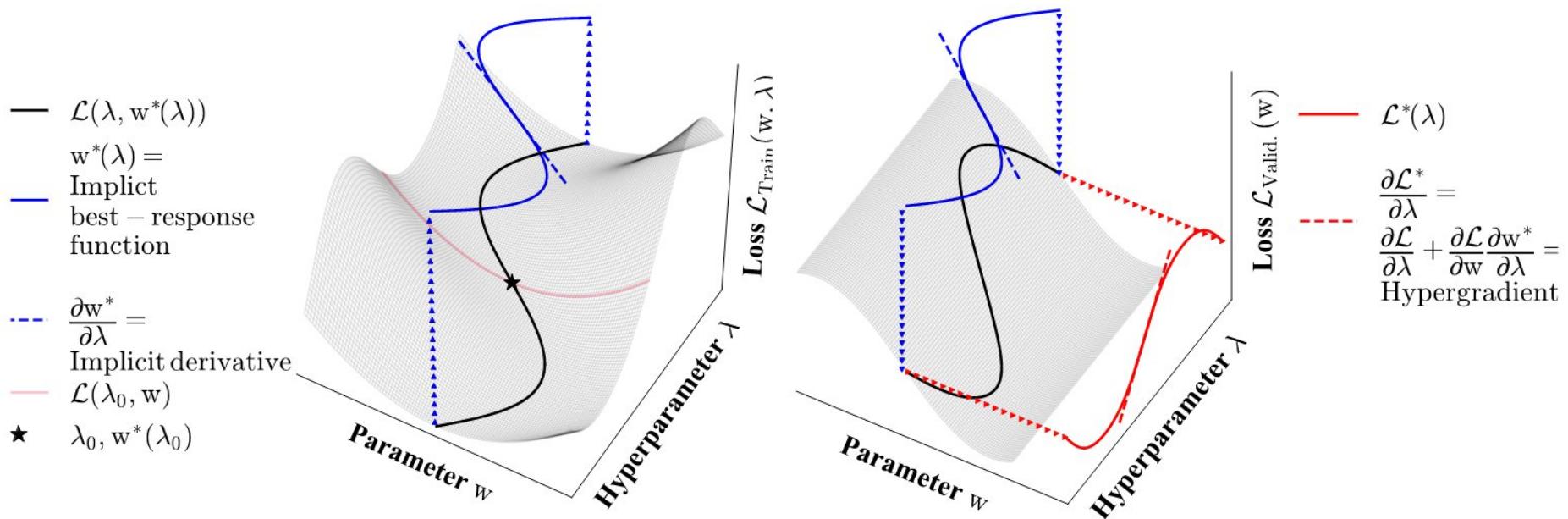
- Tree Parzen Estimator / Gaussian Process
- Sequential but allows to quickly find the global optimum.
- Proposes a new set of hyper parameters based on the scores obtained by the previous ones tested.





- Bio-inspired
- Can have fatal mutation
- **Genetic Algorithm (GA)**
- **Genetic Programming (GP)**
- **Evolution Strategy (ES)**
- **Particle Swarm Optimization (PSO)**
- **Estimation of Distribution Algorithms (EDA)**

Gradient-based optimization



Optimizing Millions of Hyperparameters by Implicit Differentiation
(<https://arxiv.org/pdf/1911.02590.pdf>)

- High dimensionality
- Bi-level optimisation

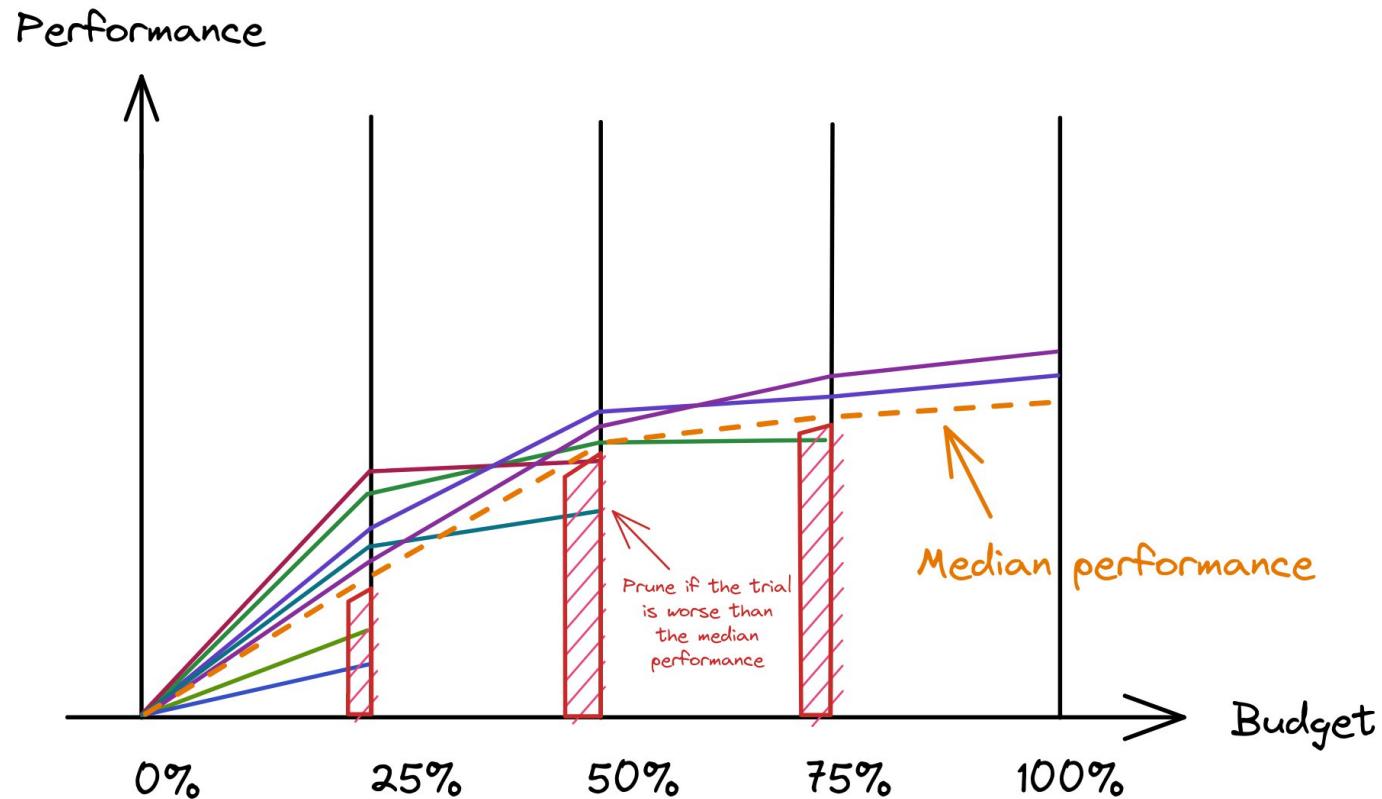
Schedulers Algorithms / Pruners

Early Stopping ◀

SHA/ASHA ◀

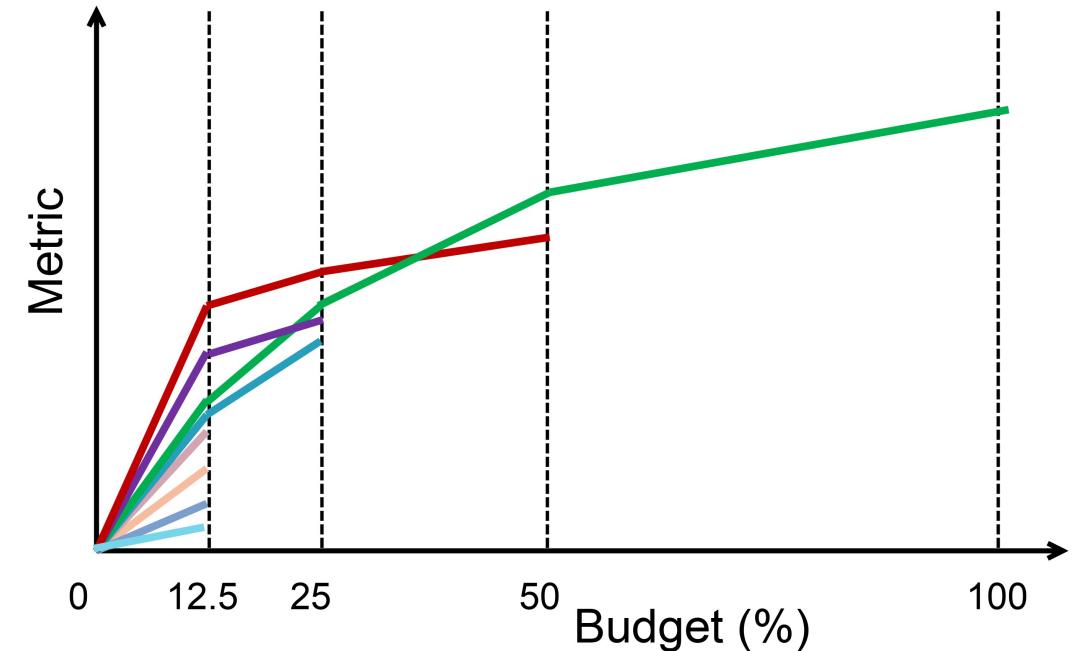
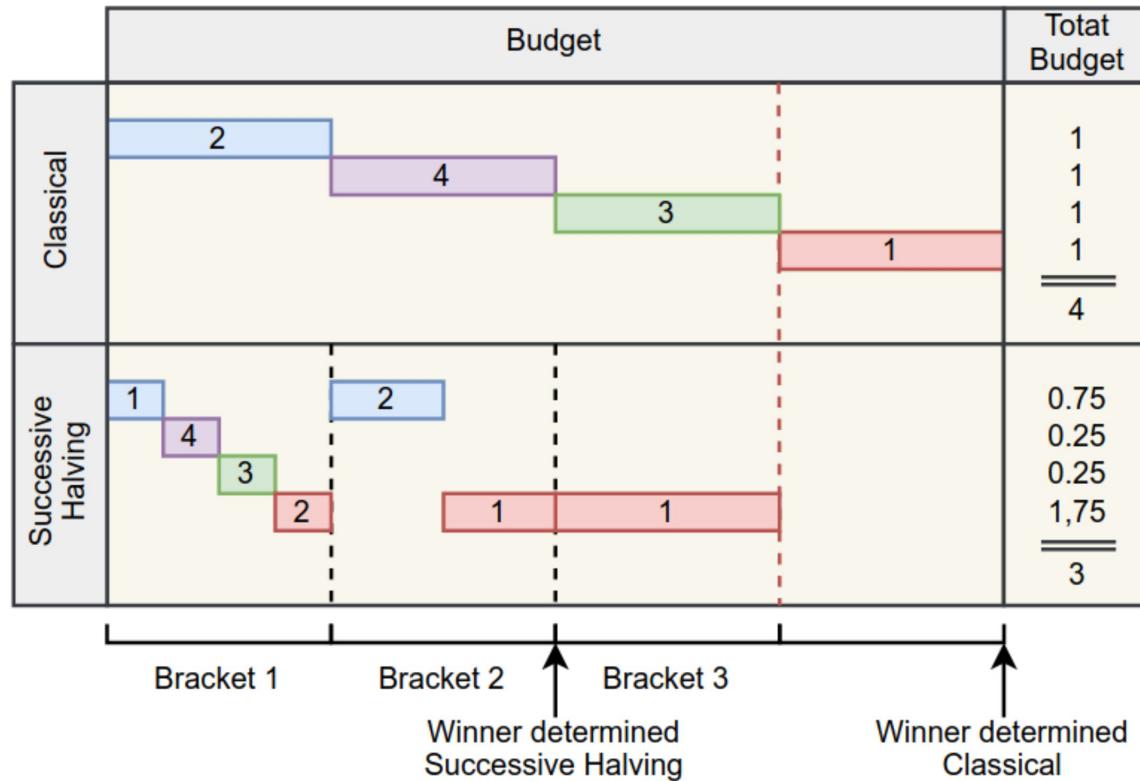
HyperBand ◀

Early Stopping

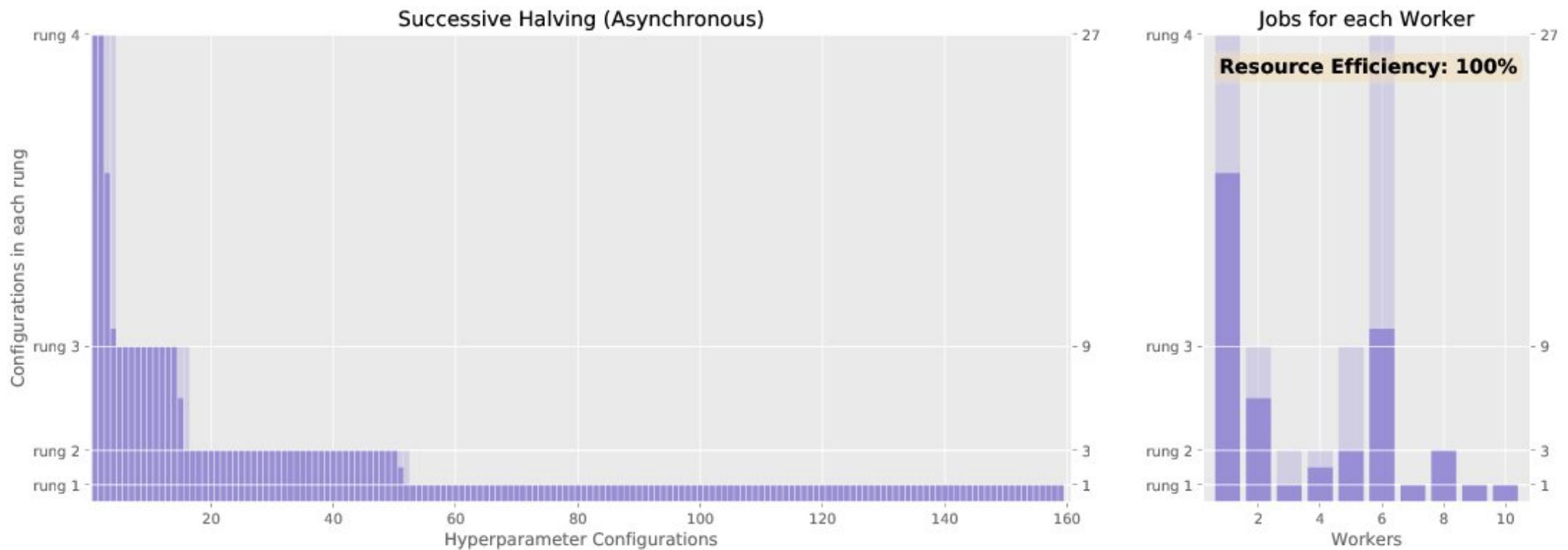


- Easy to implement
- Save resources & make automatic selection
- Can be with acc%, time%, rank%, etc

SHA : Successive Halving Algorithm



- For sequential trials
- Works well with small or medium model -> Trials must be fast !



Hyperband

Algorithm 1: HYPERBAND algorithm for hyperparameter optimization.

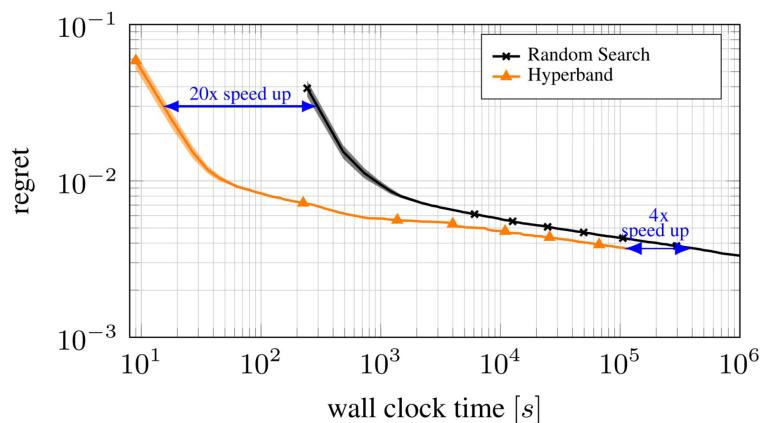
```

input      :  $R, \eta$  (default  $\eta = 3$ )
initialization:  $s_{\max} = \lfloor \log_\eta(R) \rfloor, B = (s_{\max} + 1)R$ 
1 for  $s \in \{s_{\max}, s_{\max} - 1, \dots, 0\}$  do
2    $n = \lceil \frac{B}{R(s+1)} \rceil, r = R\eta^{-s}$ 
   // begin SUCCESSIVEHALVING with  $(n, r)$  inner loop
3    $T = \text{get\_hyperparameter\_configuration}(n)$ 
4   for  $i \in \{0, \dots, s\}$  do
5      $n_i = \lfloor n\eta^{-i} \rfloor$ 
6      $r_i = r\eta^i$ 
7      $L = \{\text{run\_then\_return\_val\_loss}(t, r_i) : t \in T\}$ 
8      $T = \text{top\_k}(T, L, \lfloor n_i/\eta \rfloor)$ 
9   end
10 end
11 return Configuration with the smallest intermediate loss seen so far.

```

i	$s = 4$		$s = 3$		$s = 2$		$s = 1$		$s = 0$	
	n_i	r_i								
0	81	1	27	3	9	9	6	27	5	81
1	27	3	9	9	3	27	2	81		
2	9	9	3	27	1	81				
3	3	27	1	81						
4	1	81								

Table 1: Values of n_i and r_i for the brackets of HYPERBAND when $R = 81$ and $\eta = 3$.



- Repeatedly calls SuccessiveHalving but mitigate it's drawbacks
- Limited convergence

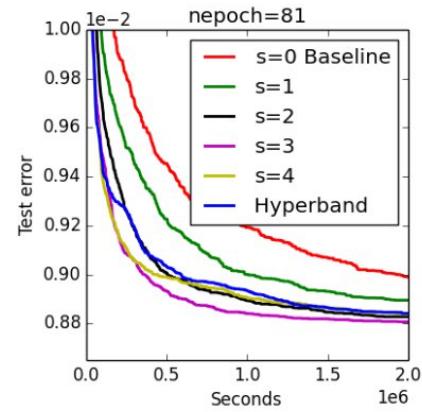


Figure 2: Performance of individual brackets s and HYPERBAND.

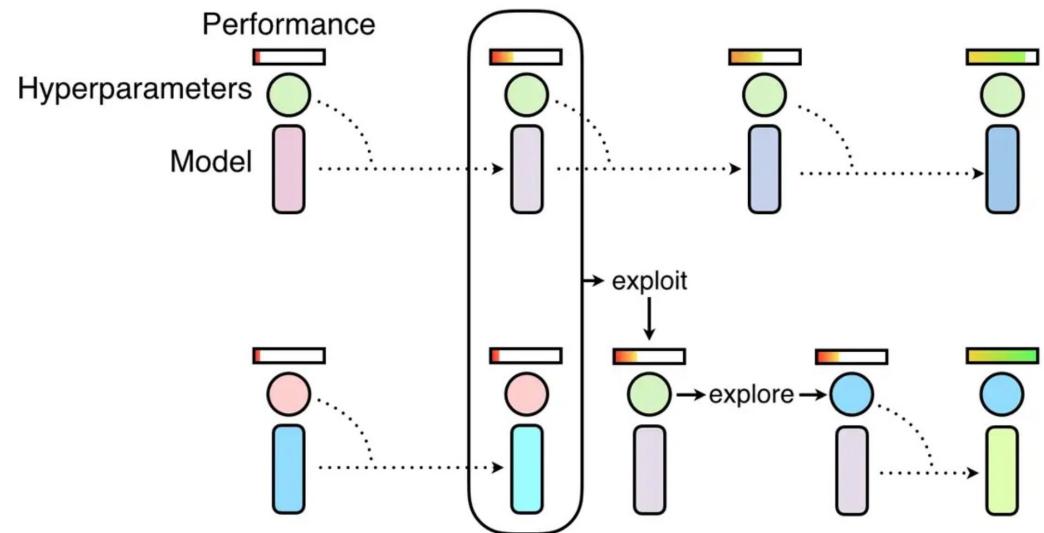
Advanced Algorithms

Hybrid time !

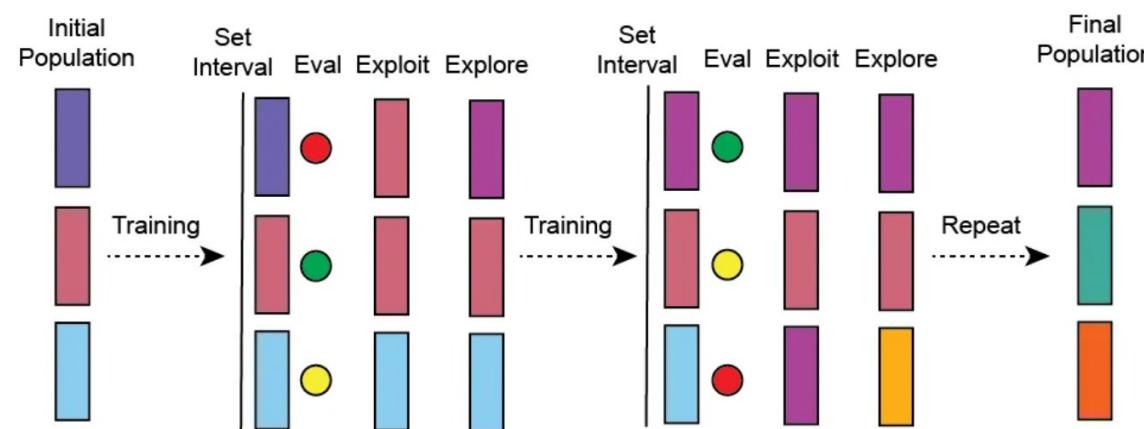
PBT ◀

BOHB, DEHB ◀

PBT : Population Based Training

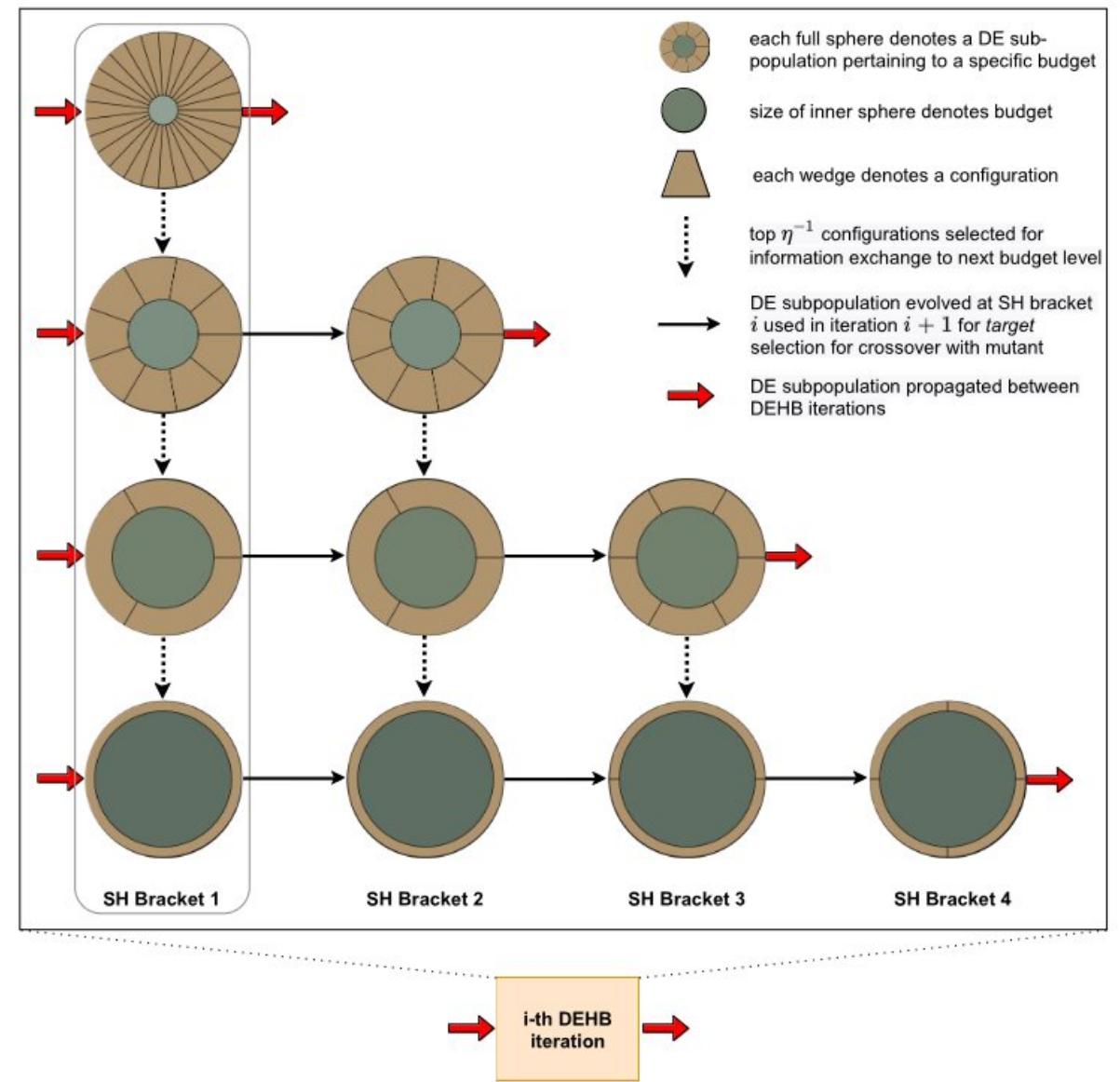
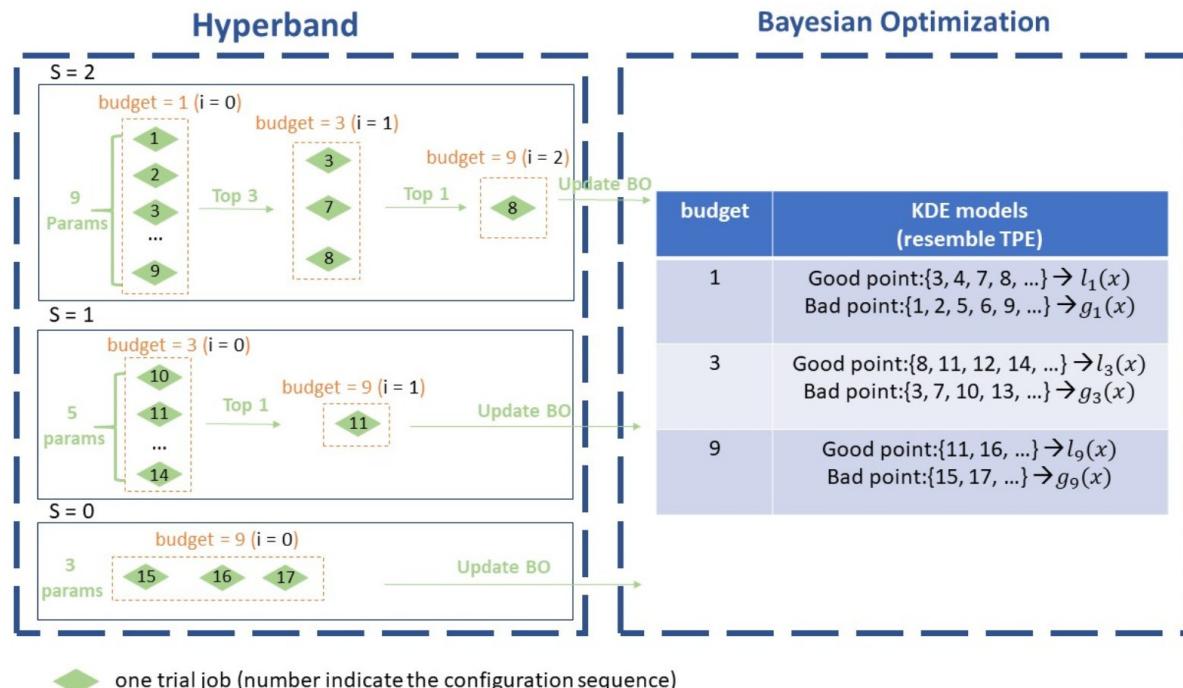


- Research and optimization of hyper parameters during training
- For large models with long and poorly parallelizable tests on a few machines.
- **Exploit** = Copy of the weights of the best model
- **Explore** = Bayesian Optimization



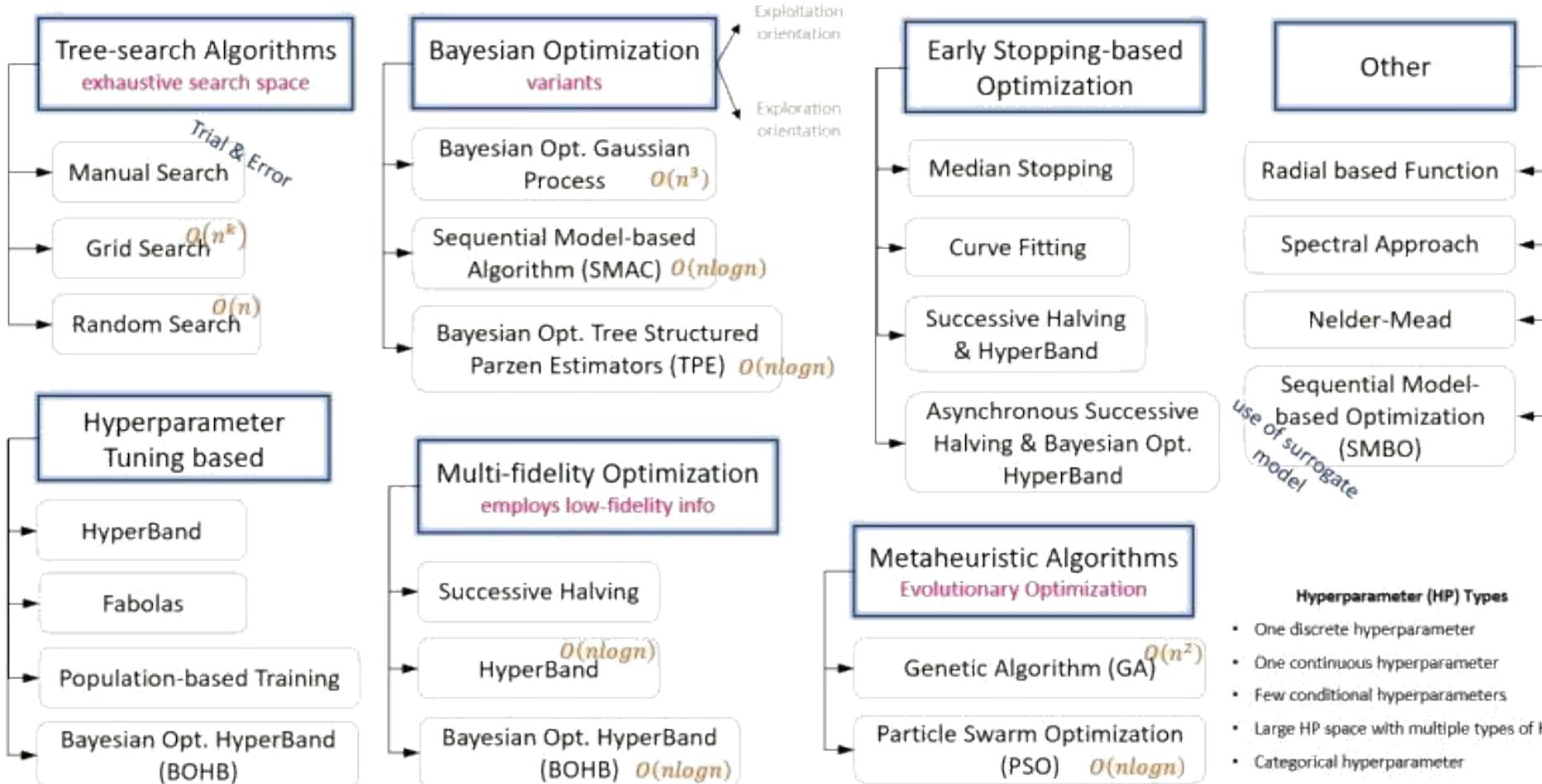
DEHB : Differential Evolution Hyperband

BOHB : Bayesian Optimization Hyperband



Selected Hyperparameter Optimization Algorithms

@Dmitry Butyrosik
The AI Vanguard
newsletter



Have the right tools

HPO frameworks ◀

Visualisation & Experiments Tracking ◀



- Based on config file
- Easy to use
- Not only used for ML/DL



OPTUNA

- Work with an objective function
- Efficient Optimization Algorithms



- Scalable HPO framework
- State of the art algorithms (PBT)
- Integrates with a wide range of additional HPO tools

 ou  Weights & Biases +

Source Control



Data Version Control



Advantages :

- allows you to save and order the results
- allows easy comparison and visualization of results
- provides all the information needed to reproduce the results

Experiments tracking

Showing 21 matching runs

				Metrics		Parameters									
	Created	Duration	Run Name	Loss/train	Loss/val	batch_size	deep	device	double	epochs	in_chan	in_dim	latent_dim	lr	residual
<input type="checkbox"/>	5 months ago	28.4min	better_lat1...	2467.9	1233.4	64	4	cuda:0	True	3	1	32	256	0.001	True
<input type="checkbox"/>	5 months ago	47.0min	better_lat1...	1840.8	525	64	3	cuda:0	True	3	1	32	128	0.001	True
<input type="checkbox"/>	5 months ago	32.2min	better_lat1...	3443.7	652.9	64	3	cuda:0	True	3	1	32	128	0.1	True
<input type="checkbox"/>	5 months ago	31.9min	better_lat1...	1680.5	3211.8	64	2	cuda:0	True	3	1	32	128	0.01	True
<input type="checkbox"/>	5 months ago	16.7min	better_lat1...	1320.8	967.3	64	3	cuda:0	True	3	1	32	128	0.01	False
<input type="checkbox"/>	5 months ago	15.2min	better_lat1...	1838.2	976.7	64	3	cuda:0	False	3	1	32	128	0.01	True
<input type="checkbox"/>	5 months ago	1.1h	better_lat2...	1239	567.6	64	3	cuda:0	True						
<input type="checkbox"/>	5 months ago	1.1h	better_lat5...	1953	1514.4	64	3	cuda:0	True						
<input type="checkbox"/>	5 months ago	18.1min	better_lat6...	1669.5	668.8	64	3	cuda:0	True						
<input type="checkbox"/>	5 months ago		better_lat1...	2368.2	3429.6	64	3	cuda:0	True						
<input type="checkbox"/>	5 months ago	45.2min	small_lat25...	1901.2	1653.1	64	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	1.1h	small_lat25...	1937.5	1039.6	64	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	29.2min	small_lat25...	2050.4	723	64	5	cuda:0	True						
<input type="checkbox"/>	5 months ago	27.2min	small_lat25...	1958.2	519.8	64	3	cuda:0	True						
<input type="checkbox"/>	5 months ago	52.5min	small_lat12...	1733.3	519.5	64	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	31.2min	small_lat25...	947.4	3987.8	64	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	1.2h	small_lat25...	1245.4	1811.4	64	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	14.1min	big-166555...	2107.9	2693.3	16	8	cuda:0	True						
<input type="checkbox"/>	5 months ago	20.0min	first-16655...	2276.2	1498.6	64	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	3.6min	first-16655...	2253.8	1435.3	32	4	cuda:0	True						
<input type="checkbox"/>	5 months ago	3.9min	first-16655...	2412.2	1407.1	32	5	cuda:0	True						

Test pretrain AutoEncoder > Comparing 6 Runs from 1 Experiment > Loss/train

Loss/train

Completed Runs 🕒
6/6

Points:

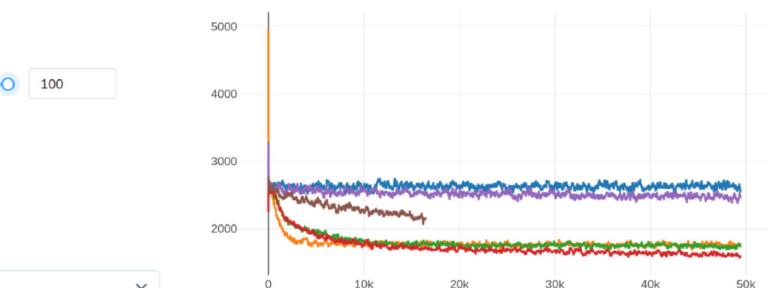
Line Smoothness 🕒
100

X-axis:
 Step
 Time (Wall)
 Time (Relative)

Y-axis:
Loss/train

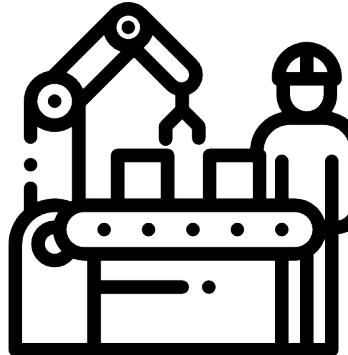
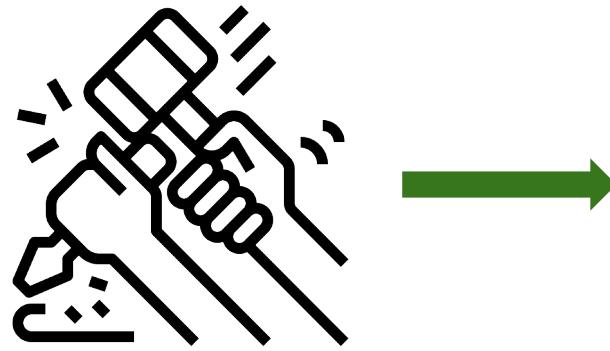
Y-axis Log Scale:

[Download CSV](#)



Loss/train

Run	Latest	Min	Max
better_lat128_deep2_lr01-1665582265	3443.7 (step=49439)	1069.9 (step=3773)	5205.4 (step=14284)
better_lat128_deep2_lr01-1665582227	1680.5 (step=49439)	806.5 (step=1210)	4966.9 (step=0)
better_lat128_deep3_lr01_nores-1665582196	1320.8 (step=49439)	870.7 (step=25165)	4310 (step=1550)
better_lat256_deep4_lr01-1665582151	1239 (step=49439)	820.8 (step=37359)	4494.6 (step=136)
small_lat256_deep5_lr01-1665575822	2050.4 (step=49439)	1045.7 (step=28182)	5172.2 (step=2325)
first-166558827	2276.2 (step=16479)	1013.7 (step=16391)	4825.8 (step=398)

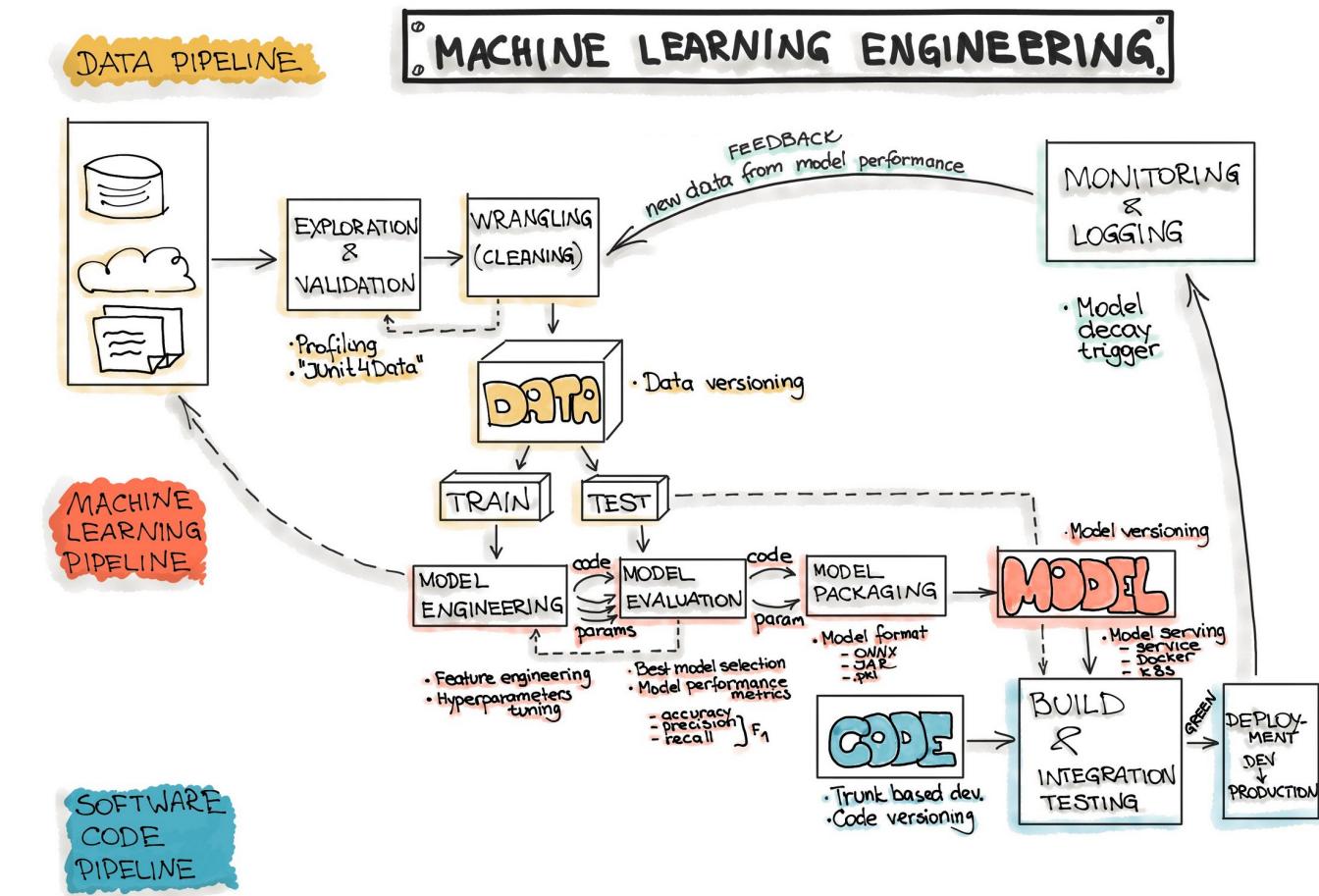


- As soon as our HPO requires a lot of resources (time, money or both) it is necessary to scale up and industrialize the experience process.
- Taking inspiration from MLOps processes and tools is a good start



Kubeflow

mlflow™



<https://ml-ops.org/content/end-to-end-ml-workflow>

- Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges (<https://wires.onlinelibrary.wiley.com/doi/pdfdirect/10.1002/widm.1484>)
- <https://www.automl.org/>
- Gradient-based Hyperparameter Optimization Over Long Horizons (<https://openreview.net/pdf?id=6x8tcREIL2W>)
- Self-Tuning networks : Bilevel Optimization of Hyperparameters using structured best-response functions (<https://openreview.net/pdf?id=r1eEG20qKQ>)
- <https://maelfabien.github.io/machinelearning/Explorium4/#>
- <https://towardsdatascience.com/a-novices-guide-to-hyperparameter-optimization-at-scale-bfb4e5047150#e813>
- Population Based Training : <https://www.deepmind.com/blog/population-based-training-of-neural-networks>