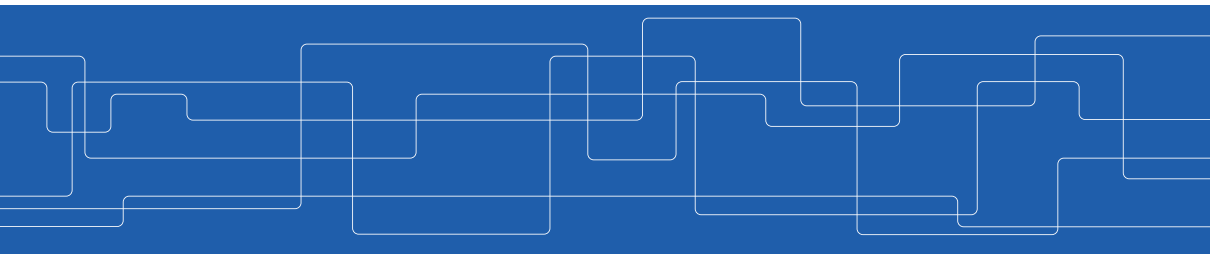




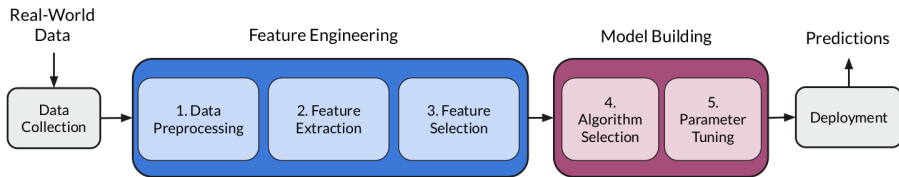
# Automated Machine Learning (AutoML)

Slides by Amir H. Payberah  
[payberah@kth.se](mailto:payberah@kth.se)



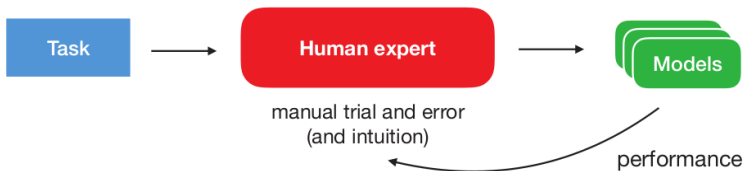
# The Machine Learning Process

- ▶ Building an **ML model** is an **iterative, complex, and time-consuming** process.
- ▶ It can take a lot of **trial and error**.

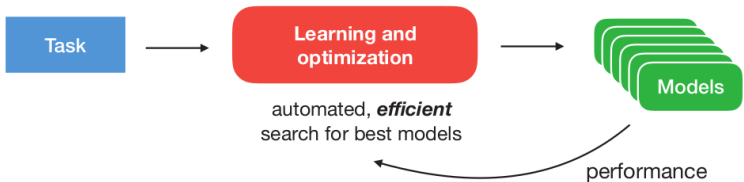


[Elshawi et al., Automated Machine Learning: State-of-The-Art and Open Challenges, 2019]

# Automated vs. Manual Machine Learning



- **AutoML**: build models in a **data-driven**, **intelligent**, and **purposeful** way.

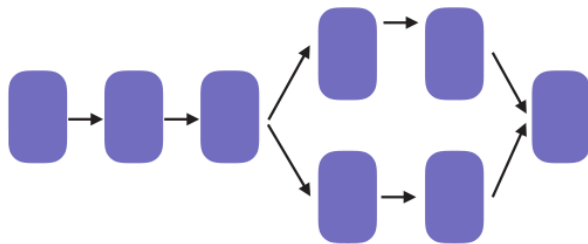


[Joaquin Vanschoren, Automatic Machine Learning - A Tutorial]



# AutoML Subproblems - Neural Architecture Search

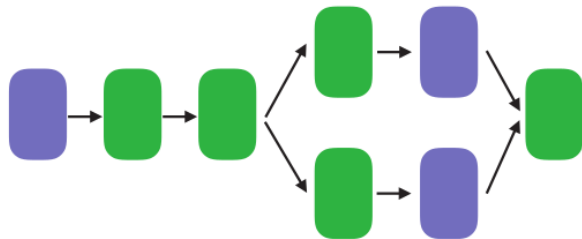
- Represent and search all pipelines or neural nets, e.g., neural layers, interconnections, etc.



[Joaquin Vanschoren, Automatic Machine Learning - A Tutorial]

# AutoML Subproblems - Hyperparameter Optimization

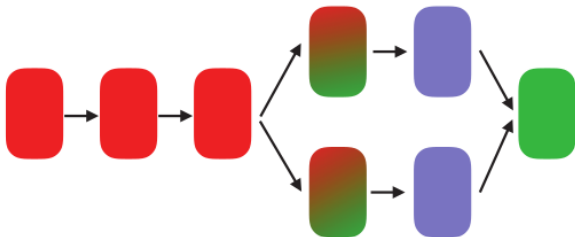
- ▶ Which **hyperparameters** are **important**? How to **optimize** them?



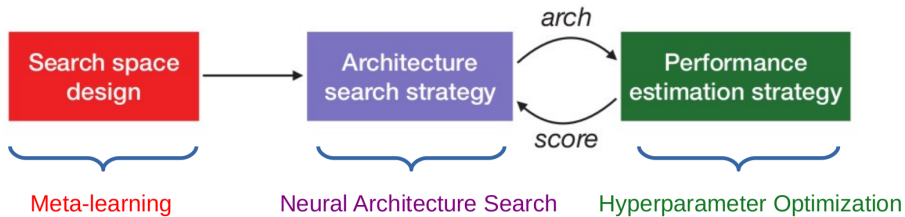
[Joaquin Vanschoren, Automatic Machine Learning - A Tutorial]

# AutoML Subproblems - Meta-learning

- ▶ How can we transfer experience from previous tasks?
- ▶ Don't start from scratch (search space is too large).



[Joaquin Vanschoren, Automatic Machine Learning - A Tutorial]







# Hyper-Parameter Optimization (HPO)



# AutoML Definition

- ▶  $A$  denotes a ML algorithm with  $m$  hyperparameters.
- ▶  $\{A_1, A_2, \dots, A_n\}$  is a set of ML algorithms.
- ▶  $\Lambda_j$  is the domain of  $j$ th hyperparameter.
- ▶  $\Lambda = \Lambda_1 \times \Lambda_2 \times \dots \times \Lambda_m$  is the overall hyperparameter configuration space.
- ▶  $\theta \in \Lambda$  is a vector of hyperparameters.
- ▶  $J(\theta, \mathbf{X}_{\text{train}}, \mathbf{X}_{\text{valid}})$  is the loss of the ML model created by  $\theta$ , trained on  $\mathbf{X}_{\text{train}}$ , and validated on  $\mathbf{X}_{\text{valid}}$ .
- ▶ Find the configuration that minimizes the expected loss on a dataset  $\mathbf{X}_{\text{train}}$ :  
$$\theta^* = \arg \min_{\theta \in \Lambda} \mathbb{E}_{(\mathbf{X}_{\text{train}}, \mathbf{X}_{\text{valid}}) \sim \mathcal{X}} J(\theta, \mathbf{X}_{\text{train}}, \mathbf{X}_{\text{valid}})$$



# Types of Hyperparameters

- ▶ **Continuous**
  - E.g., learning rate
- ▶ **Integer**
  - E.g., number of hidden units
- ▶ **Categorical**
  - E.g., choice of operator (Convolution, MaxPooling, DropOut, etc.)
  - E.g., choice of activation function (ReLU, Leaky ReLU, tanh, etc.)
- ▶ **Conditional**
  - E.g., convolution kernel size, if convolution layer is selected



# Hyper-Parameter Optimization

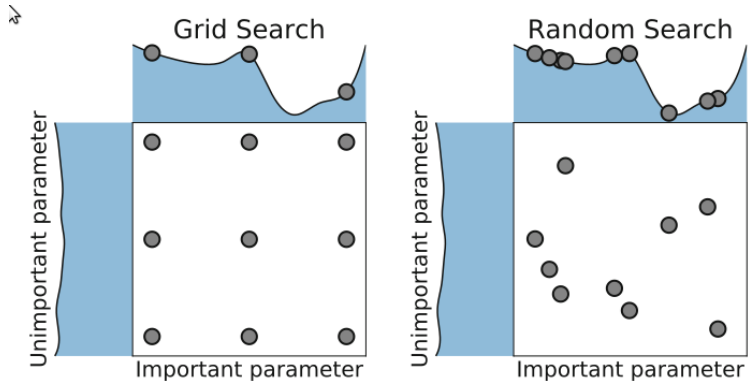
- ▶ **Black-box optimization**
  - Grid search
  - Random search
  - Population-based search
  - Bayesian optimization
  
- ▶ **Multi-fidelity optimization**
  - Modeling learning curve
  - Bandit based



# Hyper-Parameter Optimization

- ▶ **Black-box optimization**
  - Grid search
  - Random search
  - Population-based search
  - Bayesian optimization
  
- ▶ Multi-fidelity optimization
  - Modeling learning curve
  - Bandit based

# Black-box Optimization - Grid and Random Search



[Hutter et al., Automated Machine Learning, 2019]

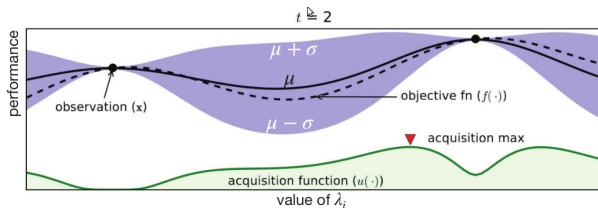


## Black-box Optimization - Population-based Search

- ▶ They maintain a **population**, i.e., a **set of configurations**.
- ▶ **Improve** this population to obtain a **new generation** of **better configurations**.
- ▶ Achieve this by applying:
  - **Local perturbations** (so-called **mutations**)
  - **Combinations** of different members (so-called **crossover**)
- ▶ E.g., genetic algorithms, evolutionary algorithms, particle swarm optimization

# Black-box Optimization - Bayesian Optimization (1/3)

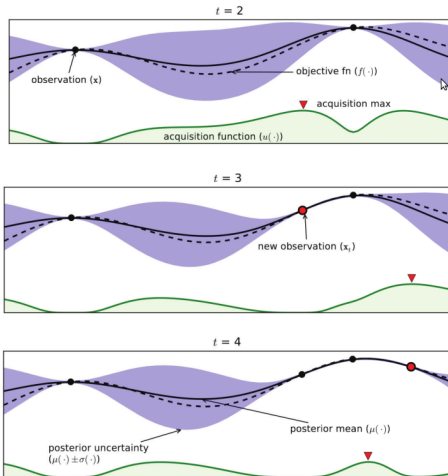
- ▶ Start with a few (random) **hyperparameter configurations**.
- ▶ Fit a **surrogate model** to **predict** other configurations.
- ▶ An **acquisition function** drives the proposition of new points to test, in an **exploration and exploitation** trade-off.
- ▶ **Sample** for the **best configuration** under that function.



[Hutter et al., Automated Machine Learning, 2019]

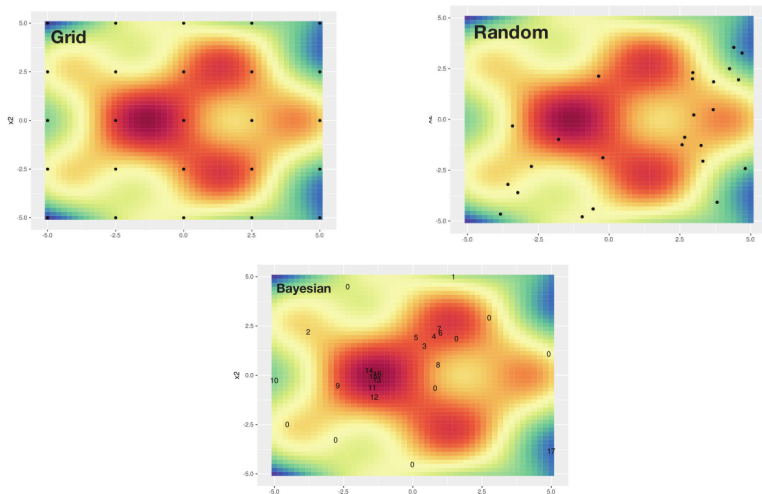


# Black-box Optimization - Bayesian Optimization (2/3)



[Hutter et al., Automated Machine Learning, 2019]

# Black-box Optimization - Bayesian Optimization (3/3)



[Hutter et al., Automated Machine Learning, 2019]



# Hyper-Parameter Optimization

- ▶ Black-box optimization
  - Grid search
  - Random search
  - Population-based search
  - Bayesian optimization
  
- ▶ Multi-fidelity optimization
  - Modeling learning curve
  - Bandit based



# Multi-fidelity Optimization

- ▶ Massive **dataset sizes** and **complex models** make **blackbox** performance evaluation **expensive**.
- ▶ Probe a hyperparameter configuration on a **small subset**.
- ▶ **Multi-fidelity** methods use **low fidelity approximations** of the **actual loss function** to minimize.
- ▶ These approximations introduce a **tradeoff** between **optimization performance** and **runtime**.



## Multi-fidelity Optimization - Modeling Learning Curves

- ▶ **Learning curve** extrapolation is used in **predicting early termination** for a particular configuration.
- ▶ Models **learning curves** during **hyper-parameter optimization**.
- ▶ Decides whether to **allocate more resources** or to **stop the training** procedure for a **particular configuration**.
- ▶ The learning process is **terminated** if the performance of the **predicted configuration** is **less than** the performance of the **best model** trained so far in the optimization process.

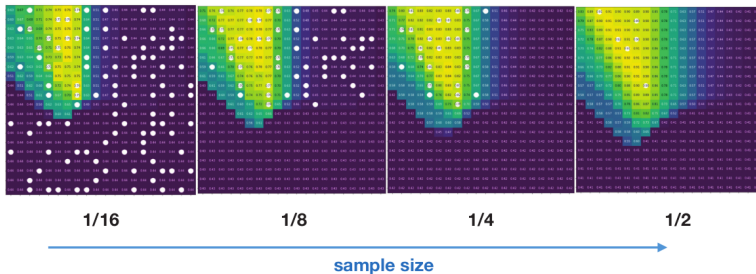


# Multi-fidelity Optimization - Bandit-Based

- ▶ Successive halving algorithm (SHA)
- ▶ HyperBand

# Multi-fidelity Optimization - SHA (1/4)

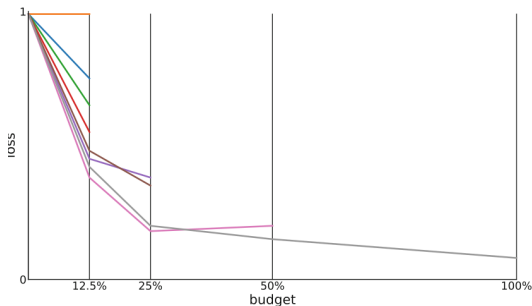
- ▶ **Train** on **small subsets**, infer which regions may be interesting to evaluate in **more depth**.
- ▶ **Randomly sample candidates** and evaluate on a small data sample.
- ▶ E.g., retrain the **50% best candidates** on twice the data.



[Hutter et al., Automated Machine Learning, 2019]

## Multi-fidelity Optimization - SHA (2/4)

- ▶ Successive halving for **eight algorithms/configurations**.
- ▶ After evaluating all algorithms on 1/8 of the total budget, **half of them are dropped** and the budget given to the remaining algorithms is **doubled**.



[Hutter et al., Automated Machine Learning, 2019]



## Multi-fidelity Optimization - SHA (3/4)

### SUCCESSIVEHALVING (Finite horizon)

**input:** Budget  $B$ , and  $n$  arms where  $\ell_{i,k}$  denotes the  $k$ th loss from the  $i$ th arm, maximum size  $R$ ,  $\eta \geq 2$  ( $\eta = 3$  by default).

**Initialize:**  $S_0 = [n]$ ,  $s = \min\{t \in \mathbb{N} : nR(t+1)\eta^{-t} \leq B, t \leq \log_\eta(\min\{R, n\})\}$ .

**For**  $k = 0, 1, \dots, s$

Set  $n_k = \lfloor n\eta^{-k} \rfloor$ ,  $r_k = \lfloor R\eta^{k-s} \rfloor$

Pull each arm in  $S_k$  for  $r_k$  times.

Keep the best  $\lfloor n\eta^{-(k+1)} \rfloor$  arms in terms of the  $r_k$ th observed loss as  $S_{k+1}$ .

**Output :**  $\hat{i}, \ell_{\hat{i},R}$  where  $\hat{i} = \arg \min_{i \in S_{s+1}} \ell_{i,R}$



## Multi-fidelity Optimization - SHA (4/4)

- ▶ Successive halving **suffers** from the **budget-vs-number** of configurations trade off.
- ▶ **Given a total budget**, the user has to decide beforehand whether:
  - to try **many configurations** and only assign a **small budget to each**, or
  - to try only **a few** and assign them a **larger budget**.
- ▶ Assigning **too small a budget** can result in **prematurely terminating good configurations**.
- ▶ Assigning **too large a budget** can result in **running poor configurations too long** and thereby wasting resources.



## Multi-fidelity Optimization - HyperBand (1/2)

- ▶ HyperBand combats **SHA problem** when selecting from **randomly sampled configurations**.
- ▶ It divides the **total budget** into **several combinations** of **number of configurations vs. budget** for each.
- ▶ Then it **calls SHA** on **each set** of random configurations.

## Multi-fidelity Optimization - HyperBand (2/2)

**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    $n = \lceil \frac{B}{R} \frac{\eta^s}{(s+1)} \rceil, \quad r = R\eta^{-s}$ 
   // 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.

```

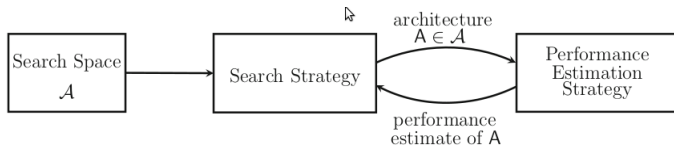
- ▶ The **inner loop** invokes **SHA** for fixed values of **n** and **r**.
- ▶ The **outer loop** iterates over different values of **n** and **r**.



# Neural Architecture Search (NAS)

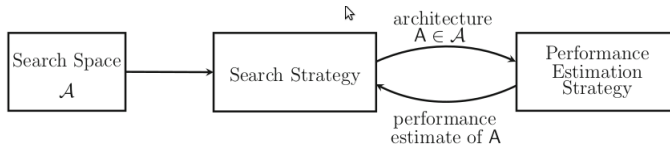
# Neural Architecture Search

- ▶ The process of **automating architecture engineering**.
- ▶ **Search space**: which architectures can be represented in principle.
- ▶ **Search strategy**: how to explore the **search space**.
- ▶ **Performance estimation**: to perform a **standard training** and **validation** of the architecture on data.



[Hutter et al., Automated Machine Learning, 2019]

# Search Space





# Search Space

- ▶ Which neural architectures a NAS approach might discover.
- ▶ Chain-structured neural network
- ▶ Multi-branch networks
- ▶ Repeated motifs



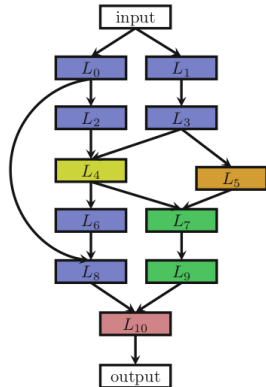
# Chain-Structured Neural Network

- ▶ A sequence of  $n$  layers.
- ▶ The  $i$ 'th layer  $L_i$  receives its **input** from layer  $i - 1$  and its **output** serves as the **input** for layer  $i + 1$ .
- ▶ **Parameters** of the search space:
  - The (maximum) **number of layers**  $n$ .
  - The **type of operation** every layer can execute, e.g., pooling, conv.
  - **Hyperparameters** associated with the **operation**, e.g., number of filters, kernel size and strides for a convolutional layer.



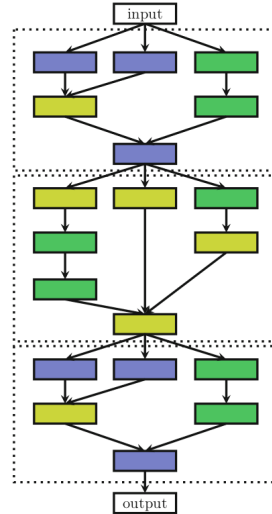
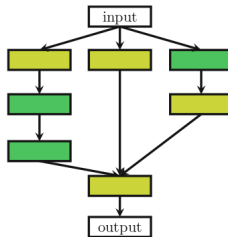
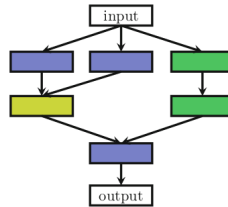
# Multi-Branch Networks

- ▶ The **input** of layer  $i$ : a function  $g_i(L_{i-1}^{\text{out}}, \dots, L_0^{\text{out}})$  of **previous layer outputs**.
- ▶ Special cases:
  - The **chain-structured** networks:  $g_i(L_{i-1}^{\text{out}}, \dots, L_0^{\text{out}}) = L_{i-1}^{\text{out}}$
  - **Residual networks**, where previous layer outputs are **summed**:  
 $g_i(L_{i-1}^{\text{out}}, \dots, L_0^{\text{out}}) = L_{i-1}^{\text{out}} + L_i^{\text{out}}, j < i$
  - **DenseNets**, where previous layer outputs are out **concatenated**:  
 $g_i(L_{i-1}^{\text{out}}, \dots, L_0^{\text{out}}) = \text{concat}(L_{i-1}^{\text{out}}, \dots, L_0^{\text{out}})$

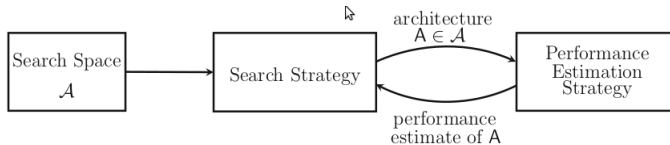


# Repeated Motifs

- ▶ **Normal cell:** preserves the dimensionality of the input.
- ▶ **Reduction cell:** reduces the spatial dimension.



# Search Strategy





# Search Strategy

- ▶ Random search
- ▶ Reinforcement learning
- ▶ Gradient-based optimization
- ▶ Bayesian optimization
- ▶ Evolutionary methods



# Random Search

- ▶ For each node in the DAG, determine what decisions must be made.
  - Choose a node as input and a corresponding operation to apply to generate the output of the node.
  - E.g., node  $i$  can take the outputs of nodes 0 to node  $i - 1$  as input.
  - E.g., choose an operation, e.g., tanh, relu, sigmoid to apply to the output of node  $i$ .
- ▶ Sample uniformly from the set of possible choices for each decision that needs to be made.
- ▶ Moving from node to node.

[Li et al., Random Search and Reproducibility for Neural Architecture Search, 2020]



# Evolutionary Methods

- ▶ Evolves a **population of models**, i.e., a set of (possibly trained) networks.
- ▶ In every **evolution step**, at least **one model from the population** is sampled and serves as a **parent** to generate **offsprings** by applying **mutations** to it.
  - E.g., **adding** or **removing** a layer, altering the **hyperparameters** of a layer, adding **skip connections**, etc.
- ▶ **After training the offsprings**, their **fitness** (e.g., performance on a validation set) is evaluated and they are **added to the population**.
- ▶ Evolutionary methods **differ** in how they **sample** parents, **update** populations, and **generate** offsprings.



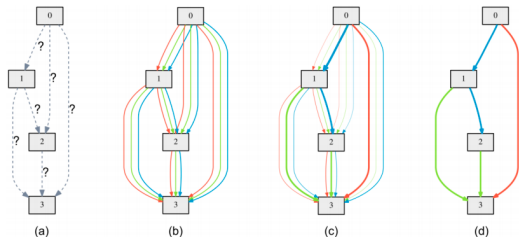
# Reinforcement Learning

- ▶ **Action**: the **generation of a neural architecture**.
- ▶ **Action space**: the **search space**.
- ▶ **Reward**: based on an estimate of the **performance** of the trained architecture on **unseen data**.
- ▶ **Policy**: different approaches.



# Gradient-based Optimization

- ▶ The **previous methods** search over a **discrete set of candidate architectures**.
- ▶ Here, it **relaxes the search space** to be **continuous**, so that the architecture can be optimized with respect to its validation set performance by **gradient descent**.
- ▶ We relax the **categorical choice** of a particular operation to a **softmax** over all possible operations.



[Liu et al., DARTS: Differentiable Architecture Search, 2019]



## Bayesian Optimization (1/3)

- ▶ Find the architecture  $\mathbf{a} \in \mathbf{A}$  that maximizes  $f(\mathbf{a})$ .
- ▶ Choose several architectures from  $\mathbf{A}$  at random and evaluating  $f(\mathbf{a})$  for each of them.
- ▶ Based on these results, iteratively choose new architectures to evaluate.
- ▶ The full algorithm:  $T$  rounds of choosing an architecture  $\mathbf{a}_i$  and computing  $f(\mathbf{a}_i)$ .
- ▶ The output is the architecture  $\mathbf{a}^*$  with the largest value of  $f(\mathbf{a}^*)$  among all those that were tried in the previous rounds.

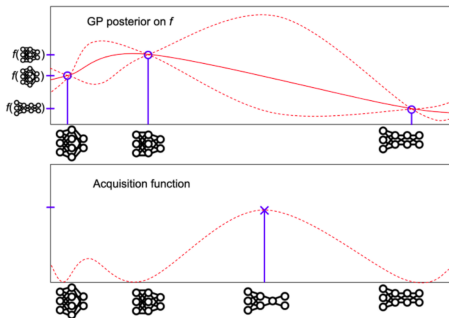


## Bayesian Optimization (2/3)

- ▶ Choose the **next architecture** in round  $i + 1$ , given  $f(a_1), \dots, f(a_i)$ .
- ▶ Assume  $f : A \rightarrow [0, 1]$  follows a **Gaussian Process (GP)**.
- ▶ Makes an assumption about the **distribution  $f(A)$** .
- ▶ The assumptions about the **mean and variance** of  $f(A)$  are **constantly being updated** as the algorithm gathers more data in the form of  $f(a_1), \dots, f(a_i)$ .
- ▶ Chooses the **architecture** with the **greatest chance** of giving a large improvement.
- ▶ The algorithm chooses  $a_{i+1} = \arg \max_{a \in A} \max(0, E[f(a) - f^*]) = \arg \max_{a \in A} E[f(a)]$ .
- ▶  $f^*$  is the **best accuracy observed so far**.

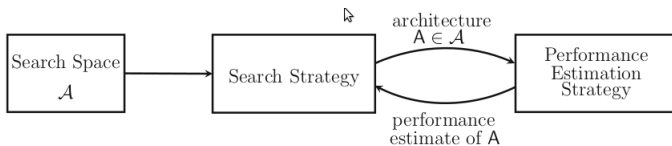
## Bayesian Optimization (3/3)

- ▶ The top graph: **three evaluations** of  $f$  (blue circles), an **estimate** of  $f$  (solid red line), and **confidence intervals** (dotted red lines).
- ▶ The bottom graph: the **expected improvement** value for each architecture. The architecture with the **largest expected improvement** is chosen (blue x).



[<https://medium.com/abacus-ai/an-introduction-to-bayesian-optimization-for-neural-architecture-search-d324830ec781>]

## Performance Estimation





# Performance Estimation

- ▶ The **search strategies** need to **estimate the performance** of a given architecture **A** they consider.
- ▶ The **simplest way** of doing this is to **train A** on **training data** and **evaluate** its performance on **validation data**.
- ▶ However, training each architecture to be evaluated **from scratch** frequently yields **computational demands** in the order of thousands of GPU days for NAS.



## Reduce the Computational Burden

- ▶ Low-fidelity approximation
- ▶ Learning curve extrapolation
- ▶ One-shot architecture



# BOHB: Robust and Efficient Hyperparameter Optimization at Scale



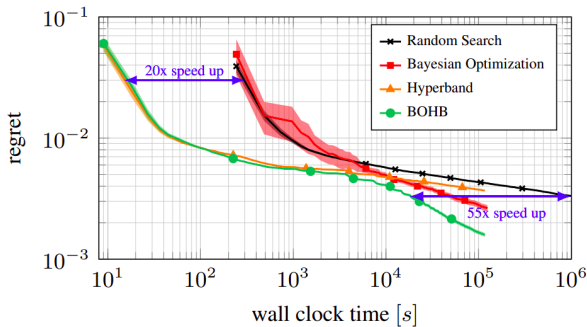


## BOHB: Bayesian Optimization and Hyperband

- ▶ Bayesian optimization (BO): for choosing the configuration to evaluate
- ▶ Hyperband (HB): for deciding how to allocate budgets

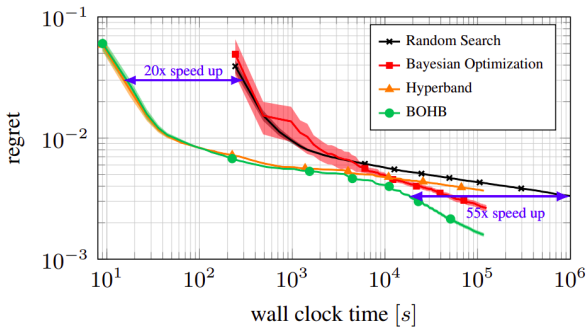
# Bayesian Optimization vs. Random Search

- ▶ **BO** advantage: much improved final performance



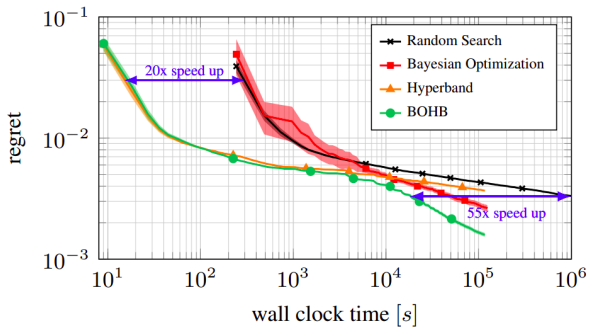
# Hyperband vs. Random Search

- ▶ **HB** advantage: much improved **anytime performance**



# Combining Bayesian Optimization and Hyperband

- ▶ Best of both worlds: strong anytime and final performance





# HBOB Algorithm

- ▶ Relies on **HB** to determine **how many configurations** to **evaluate** with which **budget**.
- ▶ Replaces the **random selection** of configurations at the **beginning of each HB iteration** by a **BO model-based search**.
- ▶ Once the **desired number of configurations** for the iteration is reached, the **SHA procedure** is carried out using these configurations.



# A System for Massively Parallel Hyperparameter Tuning



# SHA

- ▶ SHA allocates a **small budget** to **each configuration**, evaluate **all configurations** and **keep the top  $\frac{1}{\rho}$** .
- ▶ It then **increases the budget per configuration** by a factor of  $\rho$ .
- ▶ **Repeats** until the **maximum per-configuration budget** of  $R$  is reached.
- ▶ SHA requires the **number of configurations**, a **min and max resource**, a **reduction factor**, and a **minimum early-stopping rate**.



## Asynchronous SHA (ASHA)

- ▶ ASHA is a technique to **parallelize** SHA, leveraging **asynchrony** to **mitigate stragglers** and **maximize parallelism**.
- ▶ ASHA **promotes configurations** to the **next rung whenever possible**, instead of waiting for a rung to complete before proceeding to the next rung.
- ▶ If **no promotions** are possible, ASHA simply **adds a configuration to** the base rung, so that more configurations can be promoted to the upper rungs.
- ▶ Given its **asynchronous** nature it **does not require** the user to pre-specify the **number of configurations** to evaluate, but it otherwise requires the same inputs as SHA.





# DARTS: Differentiable Architecture Search



## Differentiable ARchiTecture Search (DARTS)

- ▶ Instead of searching over a **discrete set of candidate** architectures, we relax the search space to be **continuous**.
- ▶ The architecture can be **optimized** with respect to its **validation set performance** by gradient descent.



## Search Space

- ▶ It searches for a **computation cell** as the **building block** of the final architecture.
- ▶ A **cell** is a **DAG** consisting of an ordered sequence of  $N$  nodes.
- ▶ Each **node**  $x^{(i)}$  is a **latent representation** (e.g. a feature map in CNNs).
- ▶ Each **directed edge**  $(i, j)$  is associated with some **operation**  $o^{(i,j)}$  that transforms  $x^{(i)}$ .
- ▶ Each **intermediate node** is computed based on all of its **predecessors**:  
$$x^{(j)} = \sum_{i < j} o^{(i,j)}(x^i)$$

## Continuous Relaxation and Optimization

- ▶ Let  $\mathcal{O}$  be a set of **candidate operations**, where each operation represents some **function**  $o$  to be applied to  $\mathbf{x}^{(i)}$ .
- ▶ To make the **search space continuous**, it relaxes the **categorical choice** of a particular operation to a **softmax** over **all possible operations**:

$$\bar{o}^{(i,j)}(\mathbf{x}) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(\mathbf{x})$$

- ▶ The **operation mixing weights** for a pair of nodes  $(i, j)$  are parameterized by a vector  $\alpha^{(i,j)}$  of dimension  $|\mathcal{O}|$ .
- ▶ At the **end of search**, a **discrete architecture** can be obtained by **replacing** each mixed operation  $\bar{o}^{(i,j)}$  with the **most likely operation**, i.e.,  $o^{(i,j)} = \arg \max_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ .

# Summary



# Summary

- ▶ Hyperparameter optimization
  - Black-box optimization
  - Multi-fidelity optimization
  
- ▶ Neural architecture search
  - Search space
  - Search strategy
  - Performance estimation



## Reference

- ▶ Elshawi et al., Automated Machine Learning: State-of-The-Art and Open Challenges, 2019
- ▶ Falkner et al., BOHB: Robust and Efficient Hyperparameter Optimization at Scale, 2018
- ▶ Li et al., A System for Massively Parallel Hyperparameter Tuning, 2020
- ▶ Liu et al., DARTS: Differentiable Architecture Search, 2019

Questions?