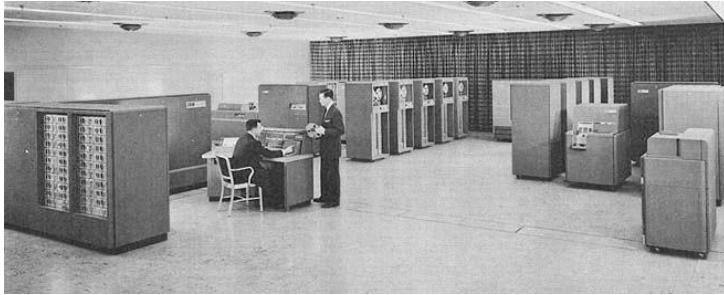


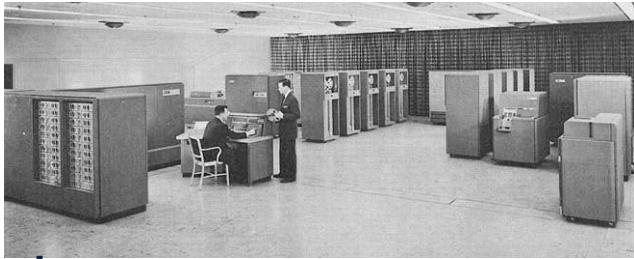
# Neural Architecture Search: The Next Frontier



Colin White  
[colin@abacus.ai](mailto:colin@abacus.ai)  
<https://crwhite.ml/>



Slides (with hyperlinks): <https://crwhite.ml/>

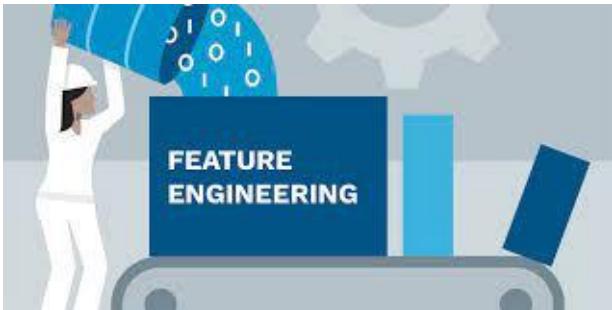


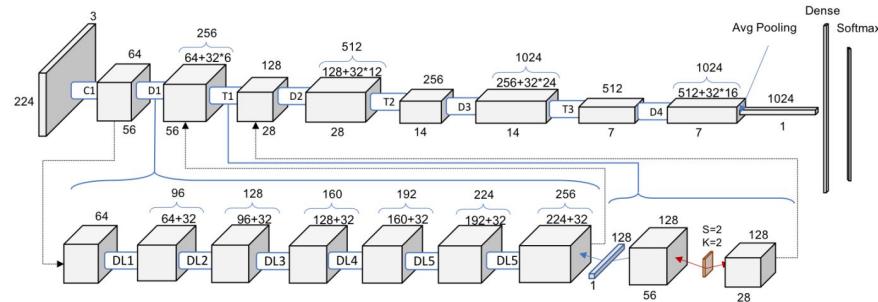
1950s

2012

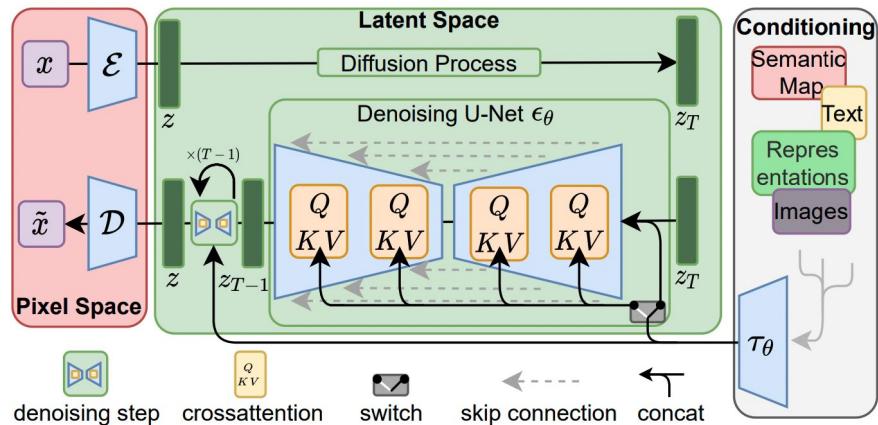
2017

2023

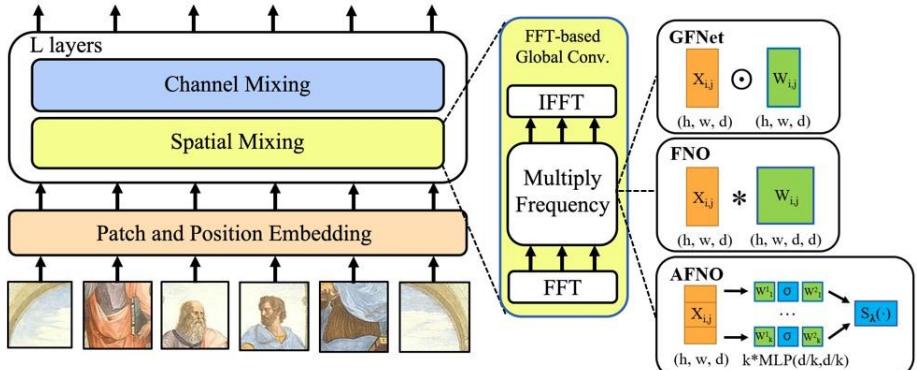




DenseNet (2016)



Stable Diffusion (2022)

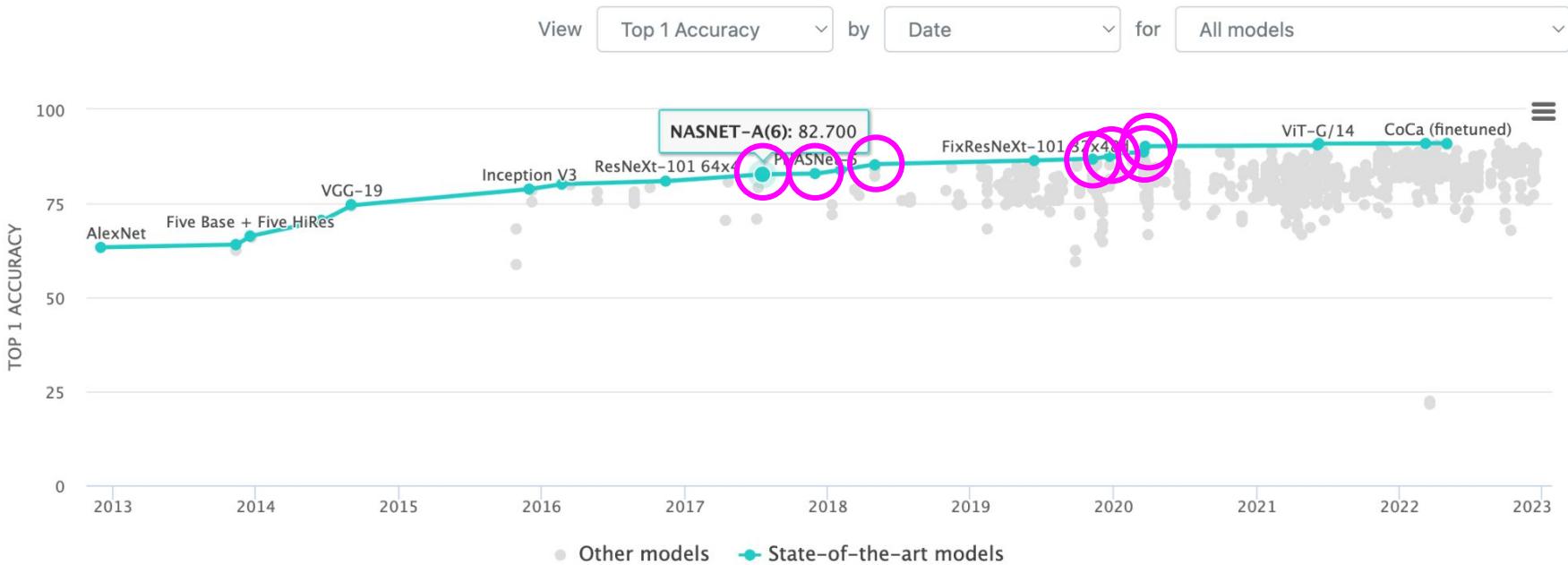


AFNO (2022)

# State of the Art on ImageNet

Leaderboard

Dataset

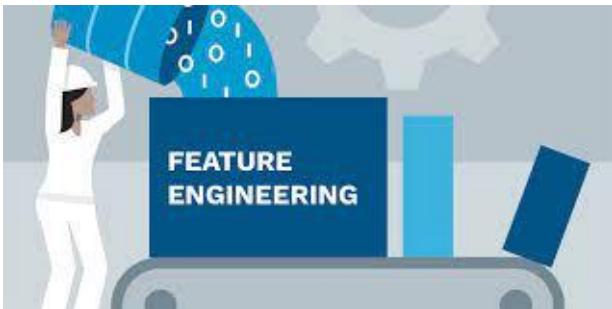


<https://paperswithcode.com/sota/image-classification-on-imagenet>



1950s

2012

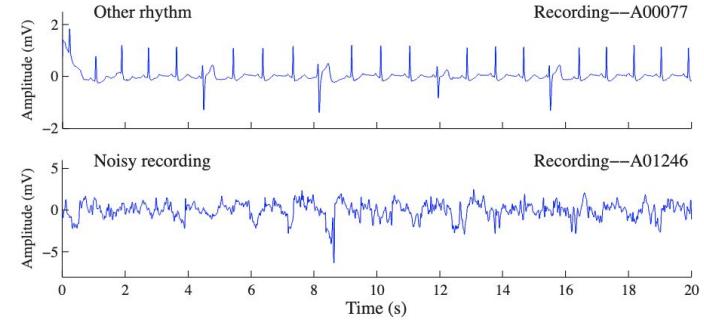
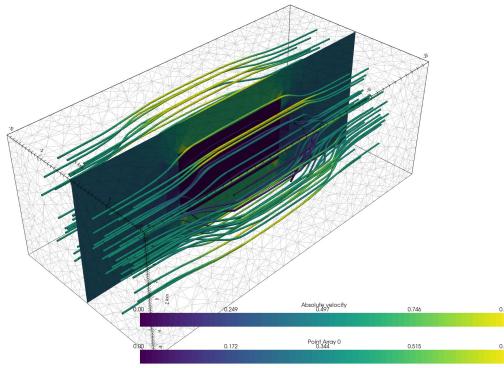
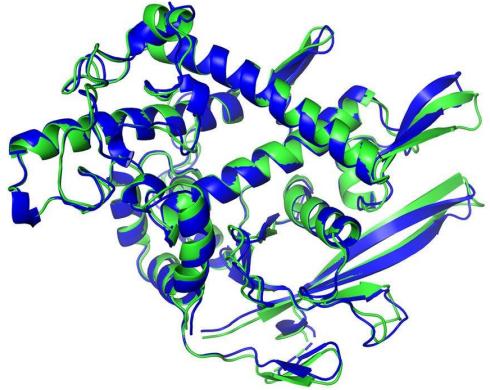


2017

2023



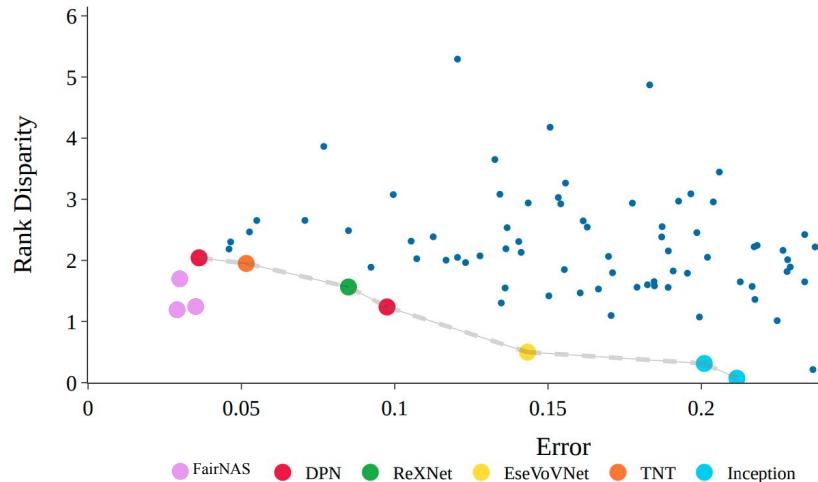
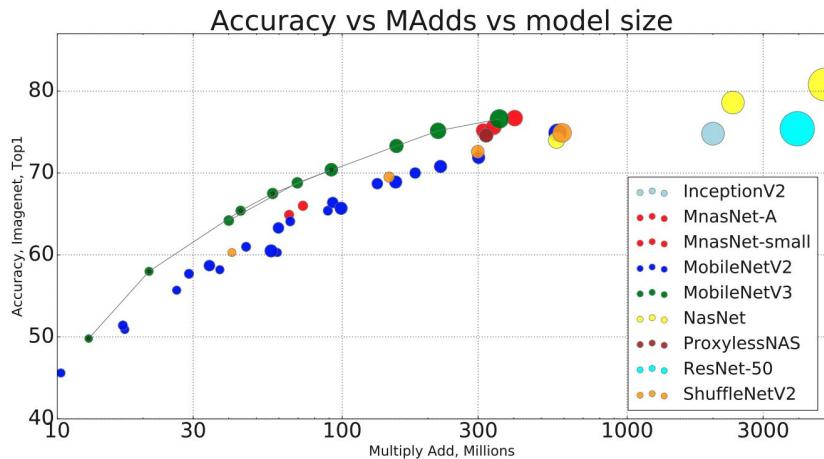
# Diverse Tasks



[TRKSST \(2022\)](#)

[TRPNJSRTNW \(2022\)](#)

# Multi-Objective



[HSCCCTWZPVLA \(2019\)](#)

[SDDWHG \(2022\)](#)



What is neural architecture search?

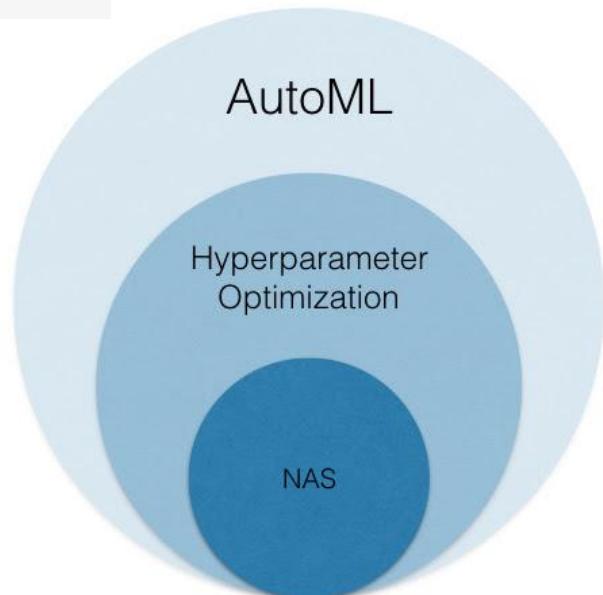


Neural architecture search (NAS) is a technique for automatically finding the best architecture for a neural network for a particular task. It involves using a search algorithm to explore the space of possible network architectures, and selecting the one that performs the best on a given dataset.

Given a search space  $A$ ,

$$\min_{a \in A} \mathcal{L}_{\text{val}}(w^*(a), a)$$

$$\text{s.t. } w^*(a) = \operatorname{argmin}_w \mathcal{L}_{\text{train}}(w, a)$$



**Automate** the design of **high-performing** models  
using **theoretical analyses, algorithm design**, and large-scale  
**empirical analyses**.

### **Neural architecture search**

WNNS (NeurIPS 2020)  
WNS (AAAI 2021)  
WNS (UAI 2021)  
Y\*W\*SH (NeurIPS 2021)  
WZRLH (NeurIPS 2021)  
M\*W\*ZKZMSYH (ICLR 2022)  
WKTSBD ICLR-Blog (2022)  
K\*W\*T\*Z\*SH (NeurIPS 2022)  
PWJNIR (arXiv 2022)  
WSSREZDH (preprint 2023)

### **Hyperparameter optimization**

BNW (COLT 2017)  
BDW (NeurIPS 2018)  
MKVDW (NeurIPS 2022)

### **Bias & explainability**

SWG (NeurIPS 2020)  
LKWN (NeurIPS 2021)  
SDDWHG (arXiv 2022)

**Enabling good scientific practices**

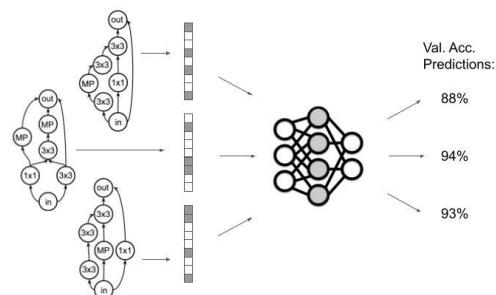
# Outline

- Introduction
- **NAS**
  - Algorithm design
  - Mitigating bias
- HPO: theoretical results
- Future directions

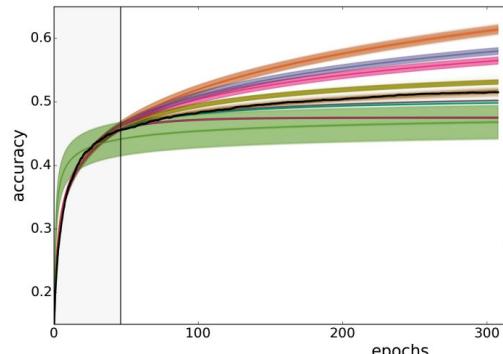


# Performance Prediction

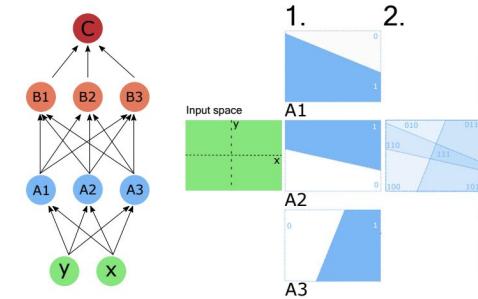
Predict the (relative) accuracy of an architecture, without fully training it



Model-based



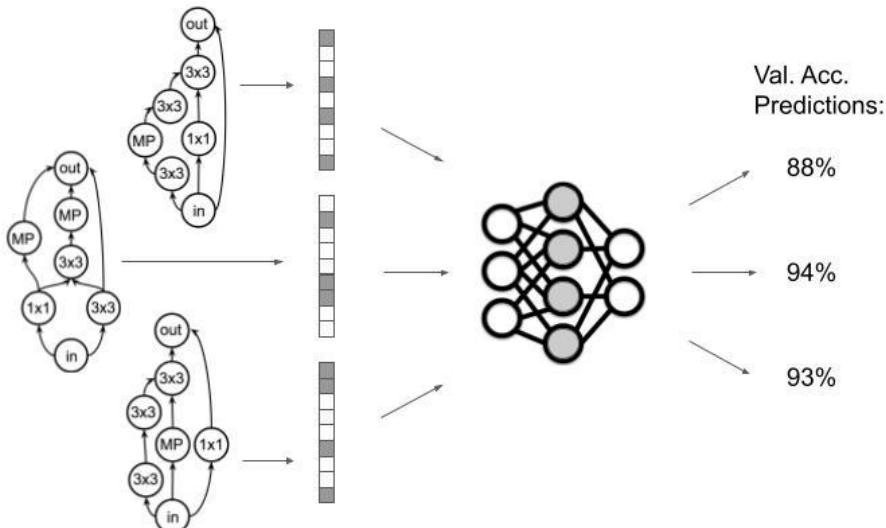
Learning curve  
extrapolation



Zero-cost proxies

# Model-Based Predictors

Train a surrogate model



- Gaussian processes
  - [KNSPX \(2018\)](#)
  - [JSH \(2018\)](#)
- Boosted trees
  - [LTWQCL \(2020\)](#)
  - [ZSZLKH \(2020\)](#)
- GNNs
  - [SPXLKZ \(2019\)](#)
  - [WLLCBK \(2019\)](#)
- Specialized encodings
  - [WNS \(2019\)](#)
  - [NZZWY \(2020\)](#)

# “BO + Neural Predictor” Components

## Algorithm 1 BANANAS

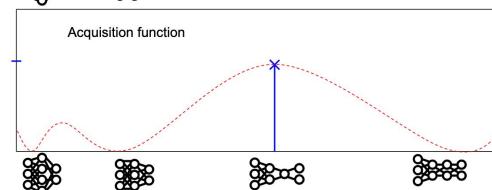
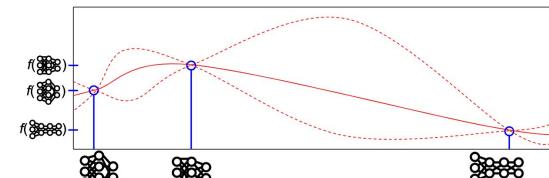
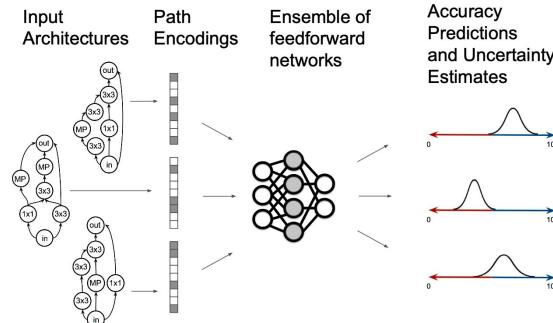
**Input:** Search space  $A$ , dataset  $D$ , parameters  $t_0, T, M, c, x$ , acquisition function  $\phi$ , function  $f(a)$  returning validation error of  $a$  after training.

1. Draw  $t_0$  architectures  $a_0, \dots, a_{t_0}$  uniformly at random from  $A$  and train them on  $D$ .
2. For  $t$  from  $t_0$  to  $T$ ,

- i. Train an ensemble of meta neural networks on  $\{(a_0, f(a_0)), \dots, (a_t, f(a_t))\}$ .
- ii. Generate a set of  $c$  candidate architectures from  $A$  by randomly mutating the  $x$  architectures  $a$  from  $\{a_0, \dots, a_t\}$  that have the lowest value of  $f(a)$ .
- iii. For each candidate architecture  $a$ , evaluate the acquisition function  $\phi(a)$ .

- iv. Denote  $a_{t+1}$  as the candidate architecture with minimum  $\phi(a)$ , and evaluate  $f(a_{t+1})$ .

**Output:**  $a^* = \operatorname{argmin}_{t=0, \dots, T} f(a_t)$ .



# “BO + Neural Predictor” Components

---

## Algorithm 1 BANANAS

---

**Input:** Search space  $A$ , dataset  $D$ , parameters  $t_0, T, M, c, x$ , acquisition function  $\phi$ , function  $f(a)$  returning validation error of  $a$  after training.

1. Draw  $t_0$  architectures  $a_0, \dots, a_{t_0}$  uniformly at random from  $A$  and train them on  $D$ .

2. For  $t$  from  $t_0$  to  $T$ ,

i. Train an ensemble of meta neural networks on  $\{(a_0, f(a_0)), \dots, (a_t, f(a_t))\}$ .

ii. Generate a set of  $c$  candidate architectures from  $A$  by randomly mutating the  $x$  architectures  $a$  from  $\{a_0, \dots, a_t\}$  that have the lowest value of  $f(a)$ .

iii. For each candidate architecture  $a$ , evaluate the acquisition function  $\phi(a)$ .

iv. Denote  $a_{t+1}$  as the candidate architecture with minimum  $\phi(a)$ , and evaluate  $f(a_{t+1})$ .

**Output:**  $a^* = \operatorname{argmin}_{t=0, \dots, T} f(a_t)$ .

---

- 
- Architecture encoding
  - Uncertainty calibration
  - Neural predictor
- architecture
- Acquisition optimization
  - Strategy
- Acquisition function

# BANANAS



## Algorithm 1 BANANAS

**Input:** Search space  $A$ , dataset  $D$ , parameters  $t_0, T, M, c, x$ , acquisition function  $\phi$ , function  $f(a)$  returning validation error of  $a$  after training.

1. Draw  $t_0$  architectures  $a_0, \dots, a_{t_0}$  uniformly at random from  $A$  and train them on  $D$ .
2. For  $t$  from  $t_0$  to  $T$ ,

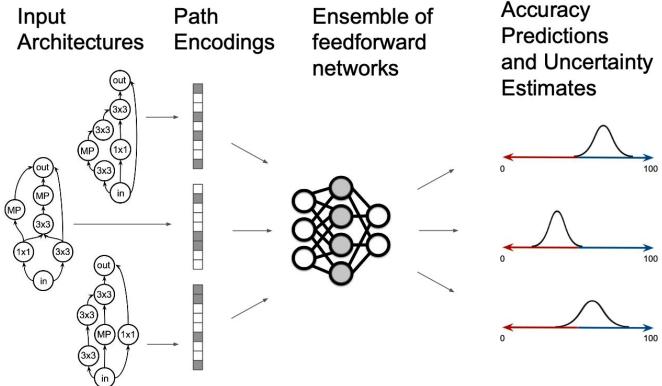
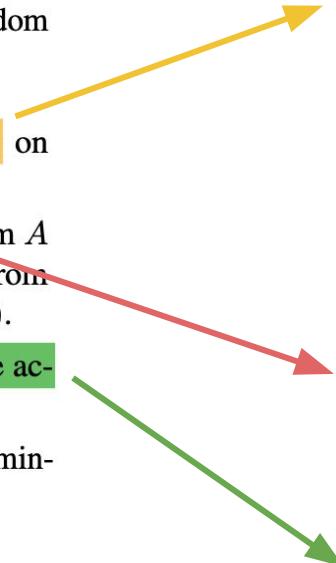
i. Train an ensemble of meta neural networks on  $\{(a_0, f(a_0)), \dots, (a_t, f(a_t))\}$ .

ii. Generate a set of  $c$  candidate architectures from  $A$  by randomly mutating the  $x$  architectures  $a$  from  $\{a_0, \dots, a_t\}$  that have the lowest value of  $f(a)$ .

iii. For each candidate architecture  $a$ , evaluate the acquisition function  $\phi(a)$ .

iv. Denote  $a_{t+1}$  as the candidate architecture with minimum  $\phi(a)$ , and evaluate  $f(a_{t+1})$ .

**Output:**  $a^* = \operatorname{argmin}_{t=0, \dots, T} f(a_t)$ .



Path encoding, ensemble

Small mutations

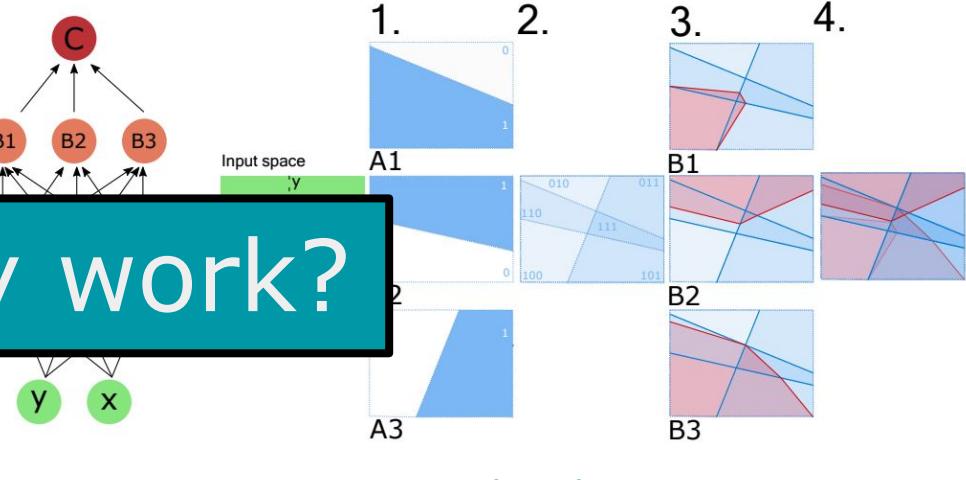
Independent Thompson Sampling

# Zero-cost proxies

|                |
|----------------|
| epe-nas [21]   |
| fisher [42]    |
| flops [25]     |
| grad-norm [1]  |
| grasp [43]     |
| l2-norm [1]    |
| jacov [23]     |
| nwot [23]      |
| params [25]    |
| plain [1]      |
| snip [14]      |
| synflow [39]   |
| zen-score [16] |

Jacobian  
Pruning-at-init  
Baseline  
Pruning-at-init  
Pruning-at-init  
Baseline  
Jacobian  
Jacobian  
Baseline  
Baseline  
Pruning-at-init  
Pruning-at-init  
Piece. Lin.

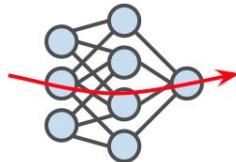
Do they work?



Compute an estimate in 5 seconds

**Examples:**  
Jacob. Cov.  
EPE-NAS  
Zen-Score

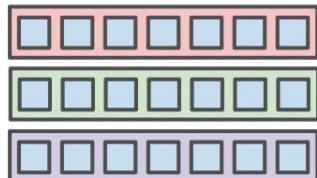
ZC Proxies



**1.5M total evaluations**

**Examples:**  
NAS-Bench-101  
TransNAS-Bench-101  
NAS-Bench-360

Benchmarks



Generalization

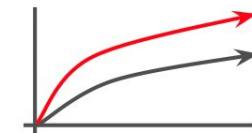
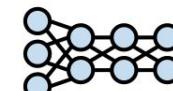


Arch. Biases

Mutual Info.

$$H(y \mid z_{i_1}, \dots, z_{i_k})$$

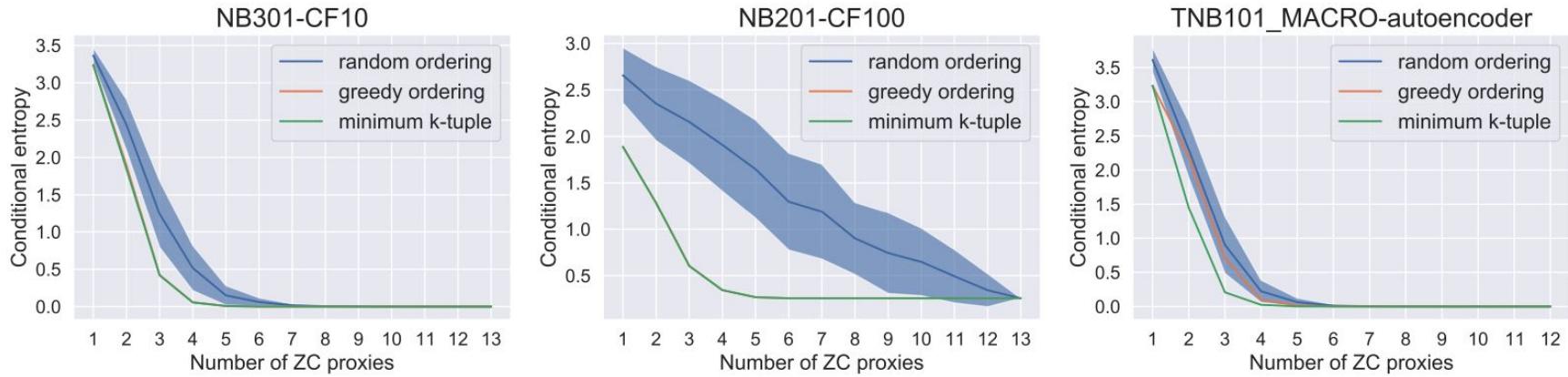
NAS Integration



[WKTSD \(ICLR-Blog 2022\)](#)

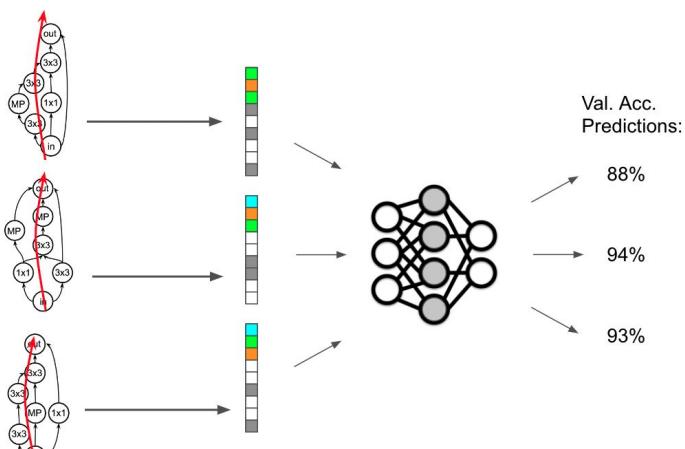
[K\\*W\\*Z\\*T\\*SH \(NeurIPS 2022\)](#)

# Complementary Information

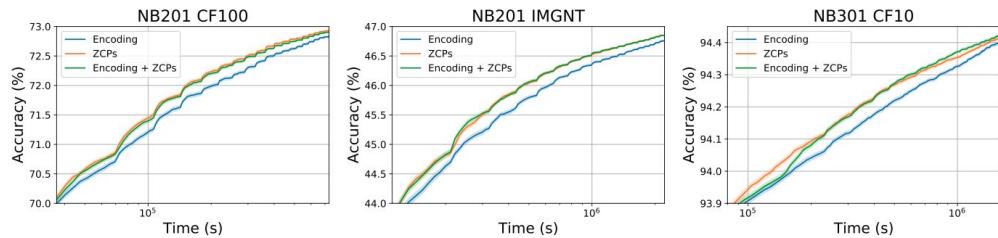


Conditional entropy  $H(y \mid z_{i_1}, \dots, z_{i_k})$  vs.  $k$

# NAS integration

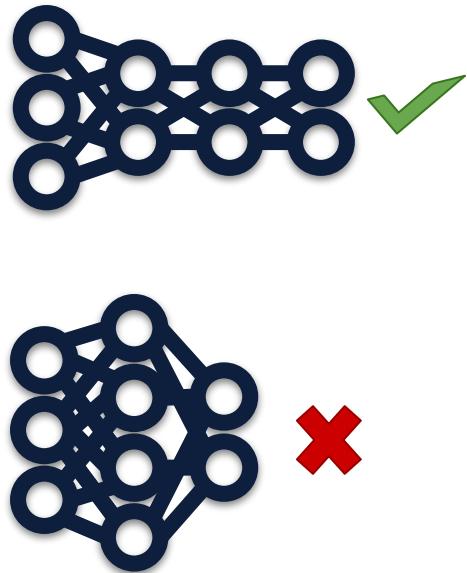


| Features<br>Benchmark | Encoding | ZC    | Both  | % Improvement (ZC) | % Improvement (Both) |
|-----------------------|----------|-------|-------|--------------------|----------------------|
| NB101-CF10            | 0.546    | 0.708 | 0.718 | 29.67              | 31.50                |
| NB201-CF10            | 0.622    | 0.905 | 0.906 | 45.50              | 45.66                |
| NB201-CF100           | 0.640    | 0.907 | 0.908 | 41.71              | 41.87                |
| NB201-IMGNNT          | 0.683    | 0.879 | 0.883 | 28.70              | 29.28                |
| NB301-CF10            | 0.314    | 0.405 | 0.465 | 28.98              | 48.09                |
| TNB101_MACRO-AUTOENC  | 0.673    | 0.831 | 0.837 | 23.48              | 24.37                |
| TNB101_MACRO-JIGSAW   | 0.809    | 0.706 | 0.809 | -12.73             | 0.00                 |
| TNB101_MACRO-NORMAL   | 0.617    | 0.710 | 0.716 | 15.07              | 16.05                |
| TNB101_MACRO-OBJECT   | 0.736    | 0.840 | 0.843 | 14.13              | 14.54                |
| TNB101_MACRO-ROOM     | 0.683    | 0.589 | 0.707 | -13.76             | 3.51                 |
| TNB101_MACRO-SCENE    | 0.832    | 0.891 | 0.899 | 7.09               | 8.05                 |
| TNB101_MACRO-SEGMENT  | 0.900    | 0.807 | 0.876 | -10.33             | -2.67                |
| TNB101_MICRO-AUTOENC  | 0.714    | 0.754 | 0.803 | 5.60               | 12.46                |
| TNB101_MICRO-JIGSAW   | 0.585    | 0.730 | 0.743 | 24.79              | 27.01                |
| TNB101_MICRO-NORMAL   | 0.657    | 0.801 | 0.809 | 21.92              | 23.14                |
| TNB101_MICRO-OBJECT   | 0.637    | 0.733 | 0.752 | 15.07              | 18.05                |
| TNB101_MICRO-ROOM     | 0.582    | 0.843 | 0.844 | 44.85              | 45.02                |
| TNB101_MICRO-SCENE    | 0.710    | 0.849 | 0.866 | 19.58              | 21.97                |
| TNB101_MICRO-SEGMENT  | 0.767    | 0.886 | 0.897 | 15.51              | 16.95                |



# Removing Biases

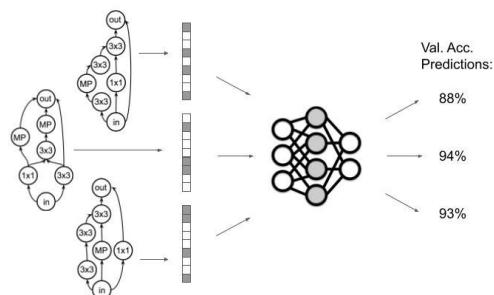
$$f'(a) = f(a) \cdot \frac{1}{b(a) + C}$$



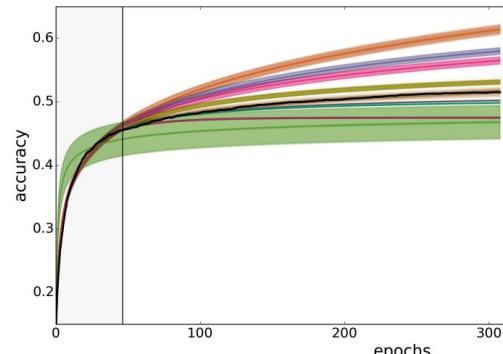
| ZC proxy | dataset     | bias metric | original bias | original perf. | new bias | new perf. | strategy    |
|----------|-------------|-------------|---------------|----------------|----------|-----------|-------------|
| l2-norm  | NB201-CF10  | conv:pool   | 0.87          | 0.42           | 0.00     | 0.10      | minimize    |
|          |             |             |               |                | 0.37     | 0.11      | equalize    |
|          |             |             |               |                | 0.70     | 0.44      | performance |
| nwot     | NB301-CF10  | conv:pool   | 0.78          | 0.49           | 0.00     | 0.03      | minimize    |
|          |             |             |               |                | 0.29     | 0.14      | equalize    |
|          |             |             |               |                | 0.78     | 0.49      | performance |
| synflow  | NB201-CF100 | cell size   | 0.57          | 0.68           | 0.01     | 0.64      | minimize    |
|          |             |             |               |                | 0.35     | 0.71      | equalize    |
|          |             |             |               |                | 0.35     | 0.71      | performance |
| synflow  | NB201-IM    | cell size   | 0.58          | 0.76           | 0.01     | 0.62      | minimize    |
|          |             |             |               |                | 0.43     | 0.76      | equalize    |
|          |             |             |               |                | 0.46     | 0.76      | performance |
| flops    | NB301-CF10  | num. skip   | -0.35         | 0.43           | -0.01    | 0.06      | minimize    |
|          |             |             |               |                | 0.12     | -0.05     | equalize    |
|          |             |             |               |                | -0.35    | 0.43      | performance |

# Performance Prediction

Predict the (relative) accuracy of an architecture, without fully training it



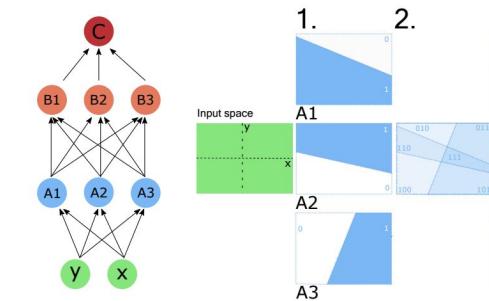
Model-based



Learning curve  
extrapolation

[WNNs \(NeurIPS 2020\)](#)  
[WNS \(AAAI 2021\)](#)

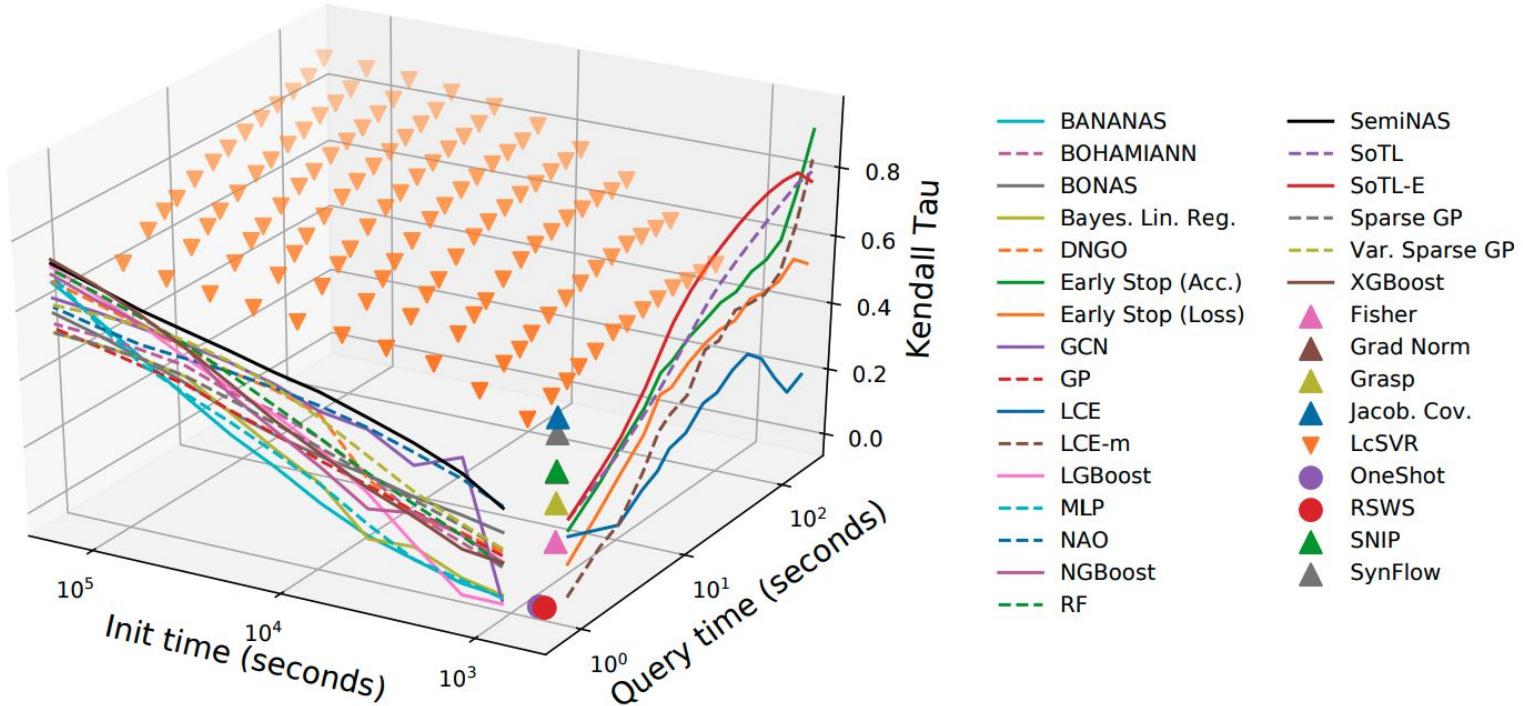
[WZRLH \(NeurIPS 2021\)](#)  
[Y\\*W\\*SH \(NeurIPS 2021\)](#)



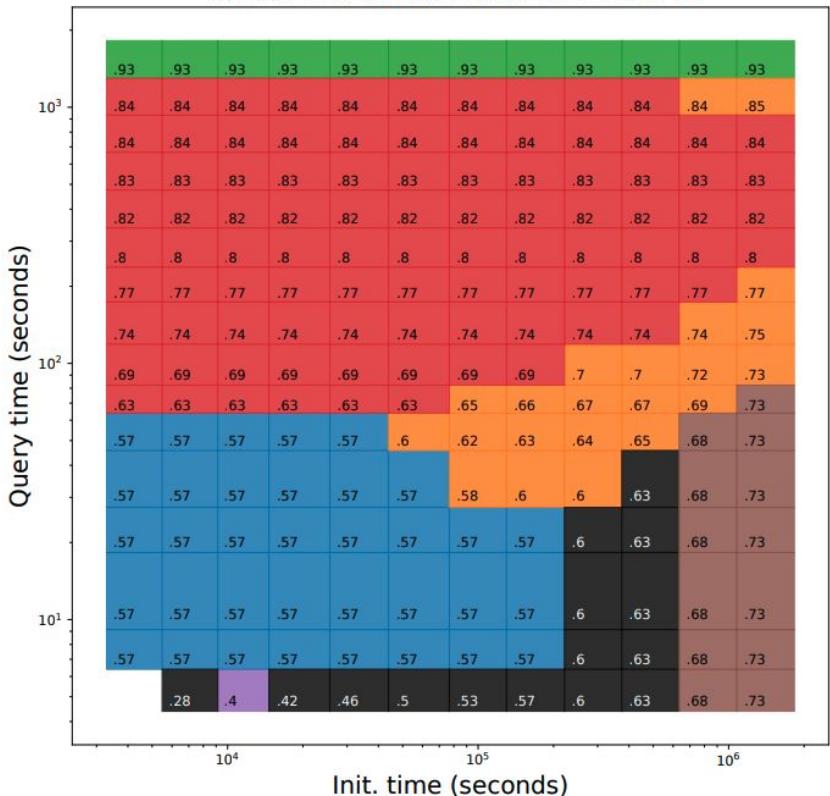
Zero-cost proxies

[WKTSD \(ICLR-Blog 2022\)](#)  
[K\\*W\\*T\\*Z\\*SH \(NeurIPS 2022\)](#)

## Kendall Tau on NAS-Bench-201 CIFAR-10

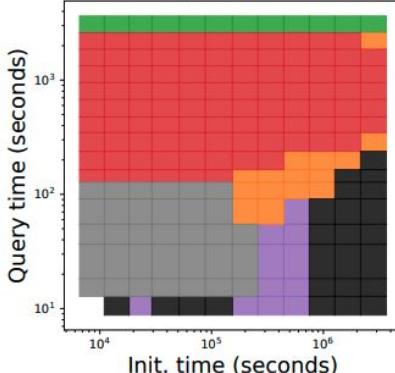


Kendall Tau on NAS-Bench-201 CIFAR-10

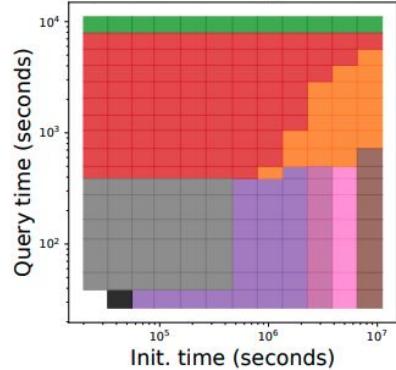


|                    |                     |               |           |           |           |
|--------------------|---------------------|---------------|-----------|-----------|-----------|
| ■ BANANAS          | ■ Early Stop (Acc.) | ■ Jacob. Cov. | ■ LcSVR   | ■ SoTL-E  | ■ SynFlow |
| ■ Bayes. Lin. Reg. | ■ GCN               | ■ LGBoost     | ■ NGBoost | ■ SemiNAS | ■ XGBoost |

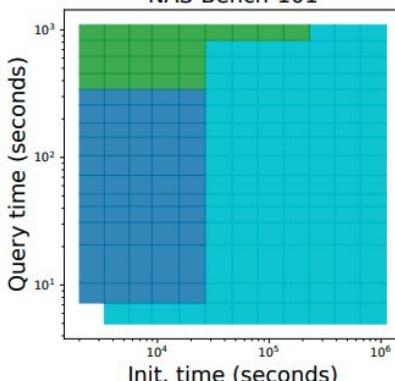
NAS-Bench-201 CIFAR-100



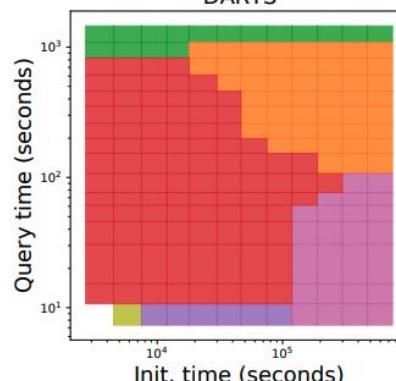
NAS-Bench-201 ImageNet16-120



NAS-Bench-101

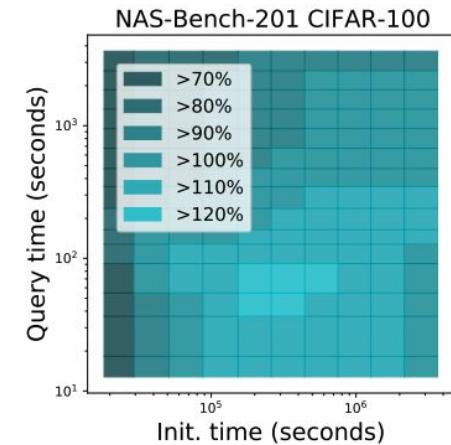
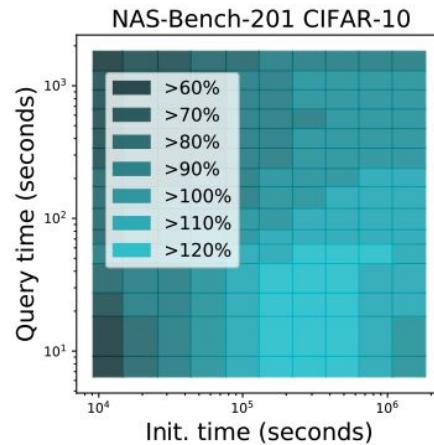
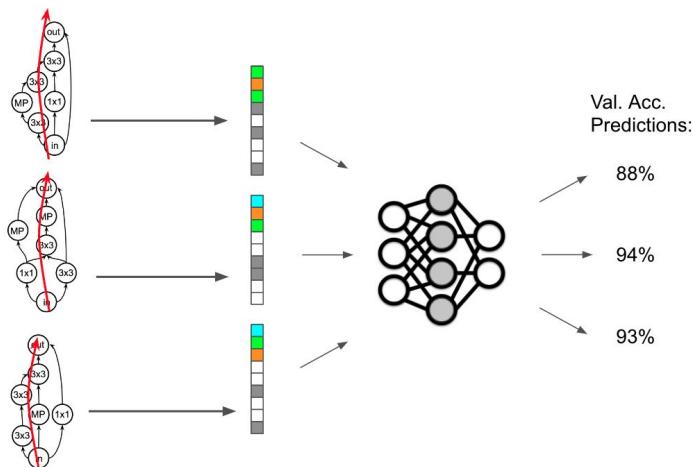


DARTS



# The Omnipotent Predictor

Combine best predictors from each family



How do we set up  
**neural architecture search**  
for a new application?

# Face Recognition

## 1. For one-to-one matching, the team saw higher rates of false positives for Asian and African American faces relative to images of Caucasians.

The differentials often ranged from a factor of 10 to 100 times, depending on the individual algorithm. False positives might present a security concern to the system owner, as they may allow access to impostors.

## 2. Among U.S.-developed algorithms, there were similar high rates of false positives in one-to-one matching

for Asians, African Americans and native groups (which include Native American, American Indian, Alaskan Indian and Pacific Islanders). The American Indian demographic had the highest rates of false positives.

## 'The Computer Got It Wrong': How Facial Recognition Led To False Arrest Of Black Man

June 24, 2020 · 8:00 AM ET



<https://www.nist.gov/news-events/news/2019/12/nist-study-evaluates-effects-race-age-sex-face-recognition-software>

## A US government study confirms most face recognition systems are racist

by Karen Hao December 20, 2019



A U.S. Customs and Border Protection officer helps a passenger navigate a facial recognition kiosk at the airport.  
DAVID J. PHILLIP/AP

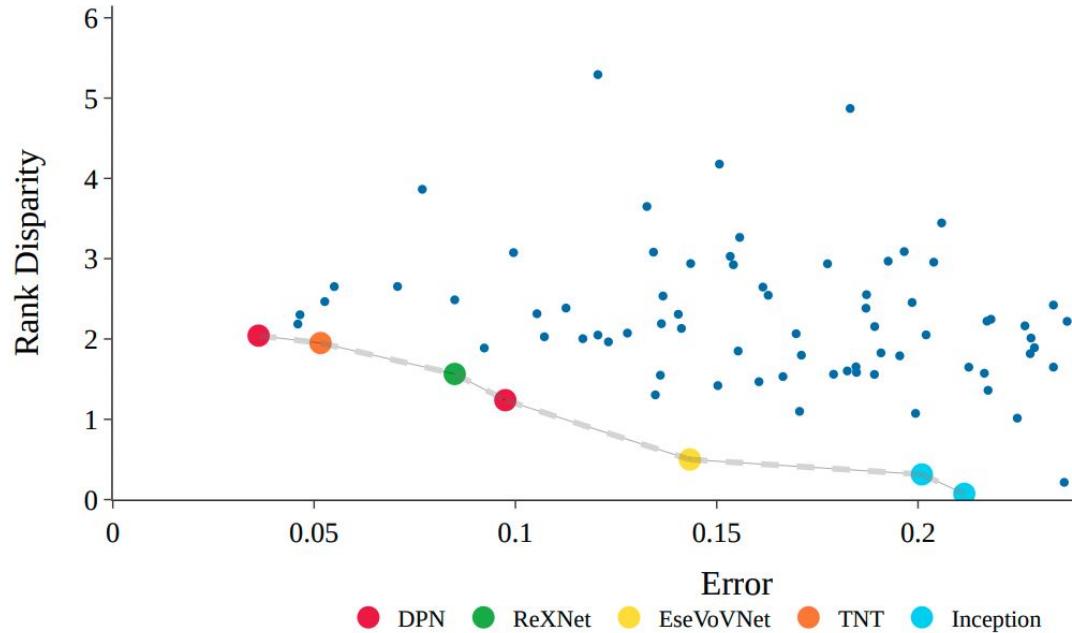
## Architectures

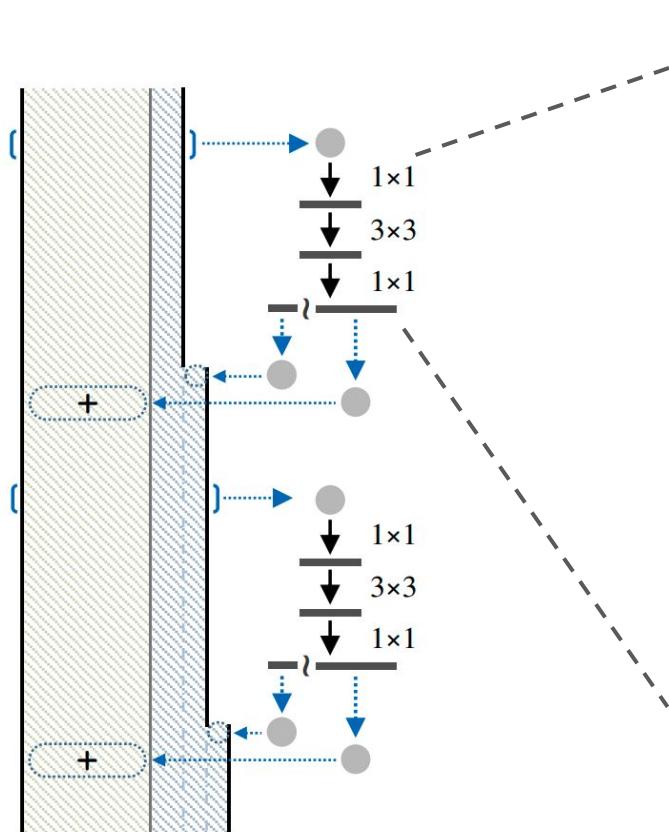
**Examples:**  
Inception  
ViTransformer  
VovNet  
MobileNetV3  
Dual Path Net

## Hyperparams

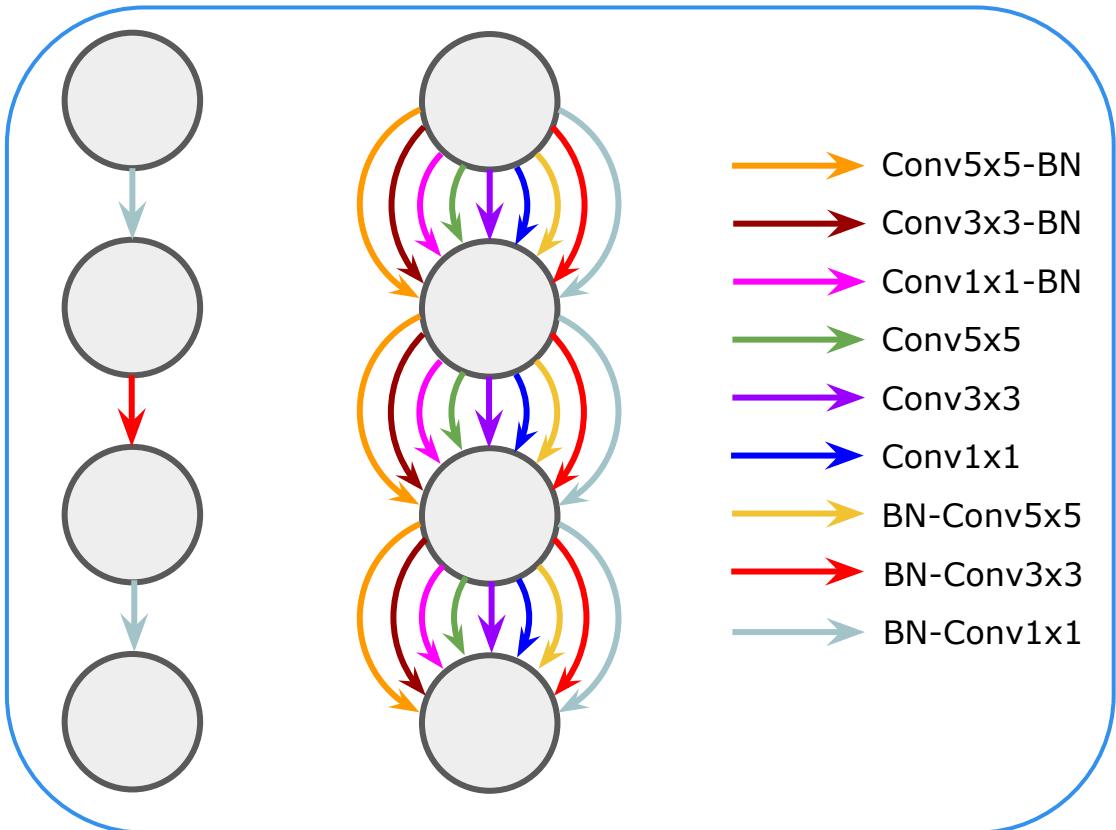


**Examples:**  
MagFace (head)  
ArcFace (head)  
CosFace (head)  
Optimizer  
Learning Rate

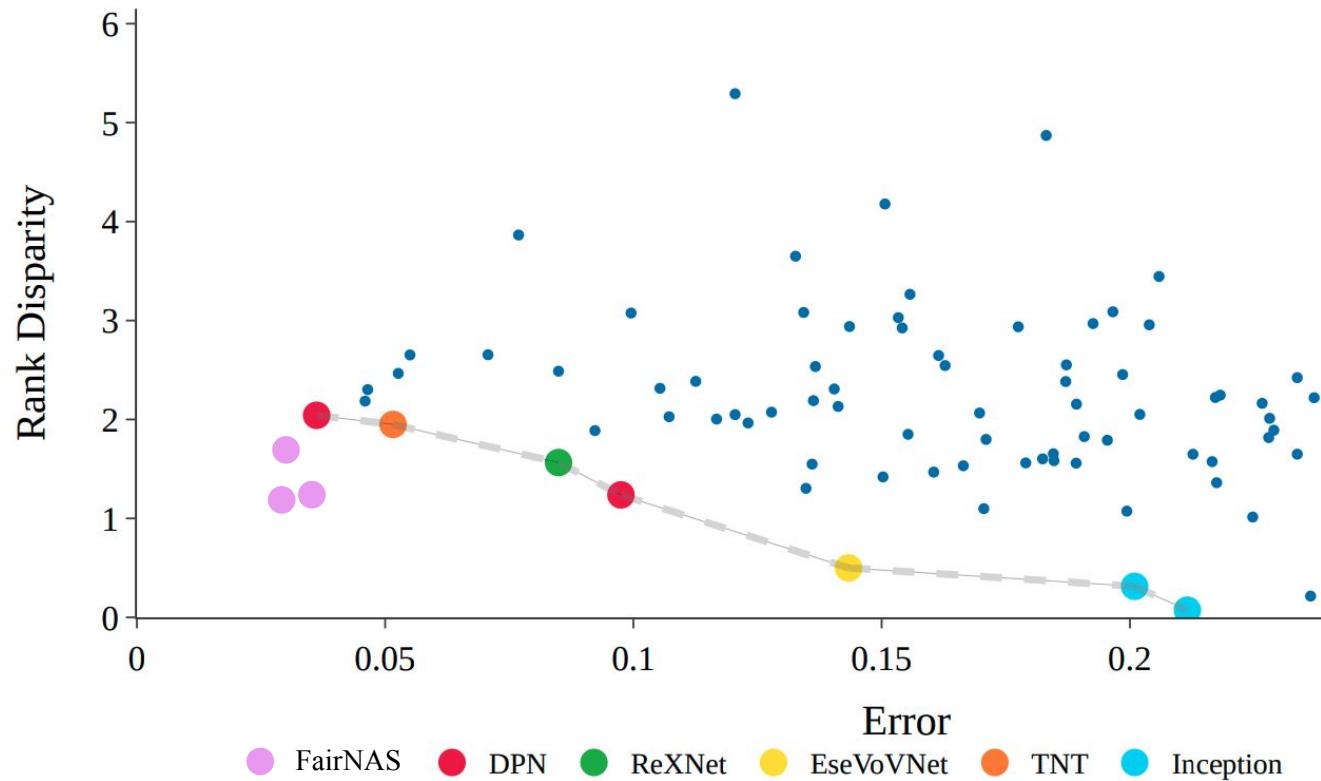




DPN (2018)



SDDWHG (arXiv 2022)



# Outline

- Introduction
- NAS: algorithm design
- **HPO: theoretical results**
- Future directions



# Data-Driven Algorithm Design



Optimization for  
**one** dataset



Optimization for a  
**distribution** of datasets

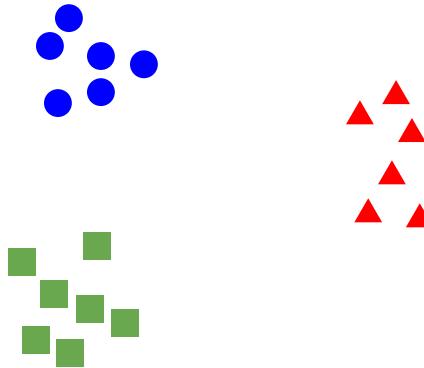
**Supervised** learning



**Unsupervised** learning  
(clustering)

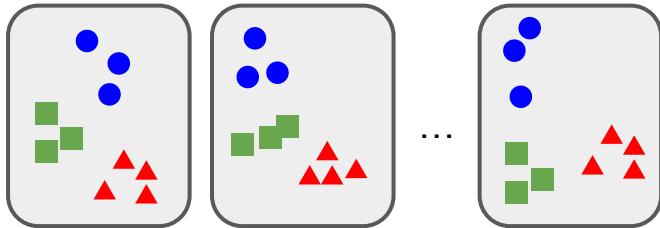
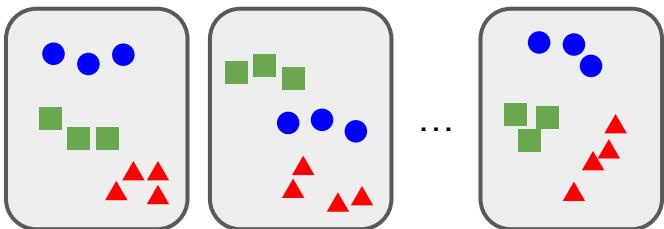
# Clustering

Input: dataset,  $k$



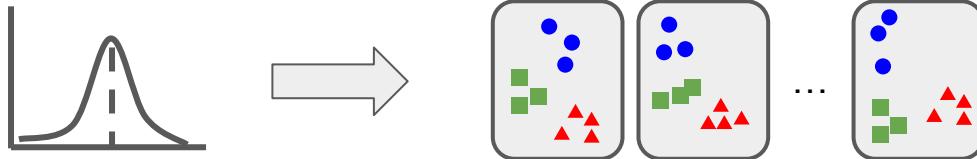
- Find  $k$  clusters
- Minimize distance to ground-truth clustering

# Data-Driven Algorithm Design



# Data-Driven Algorithm Design

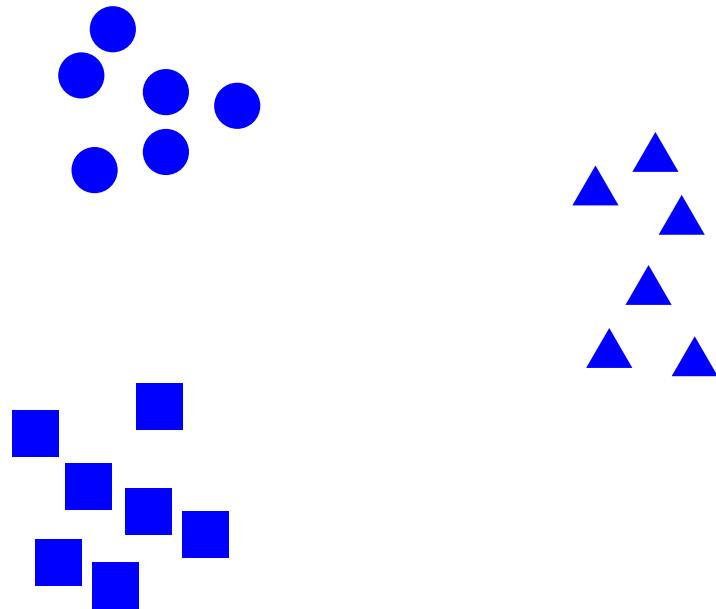
- Fix a parameterized clustering algorithm
- Receive training set of “typical” clustering instances



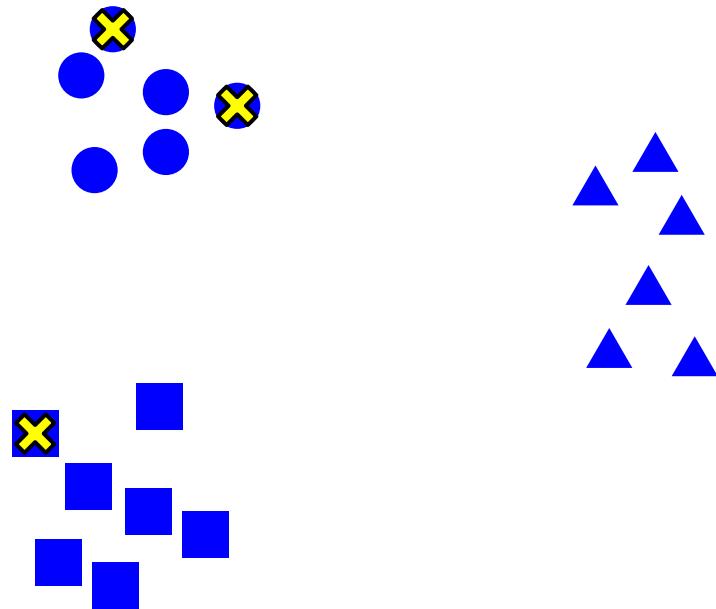
- Find hyperparameters with good average performance

How to find high-performing hyperparameters?  
Will these hyperparameters have strong future performance?

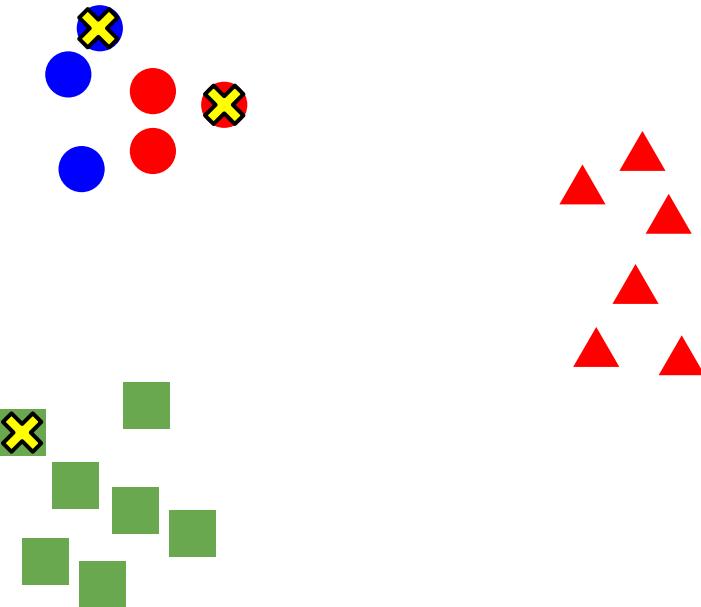
# k-means



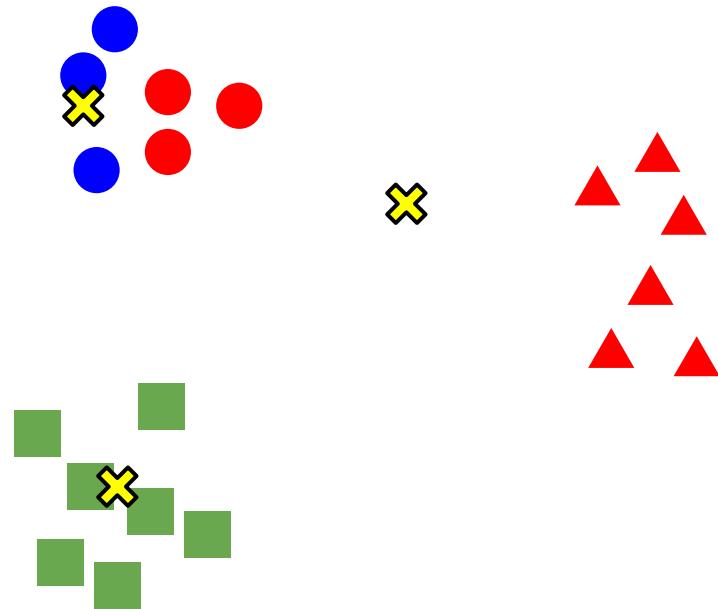
# k-means



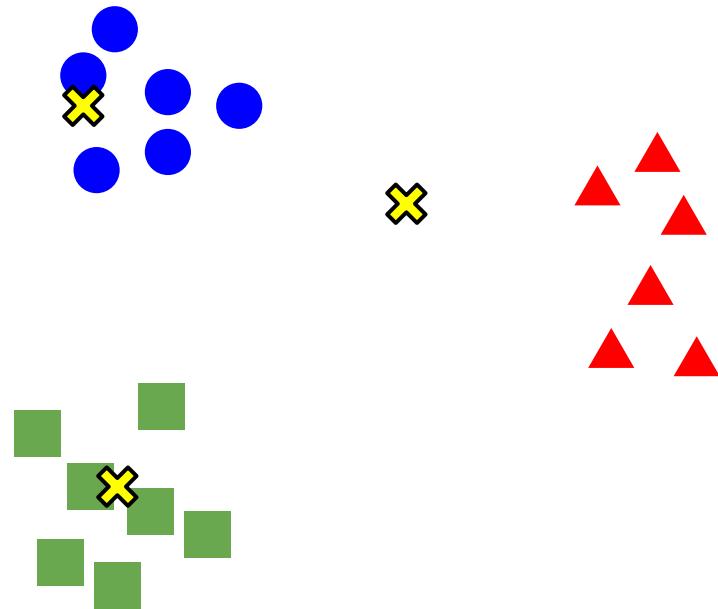
# k-means



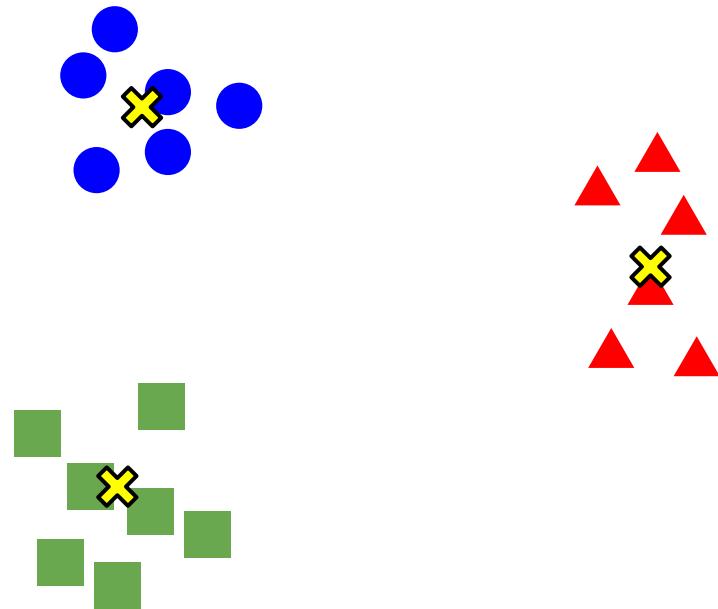
# k-means



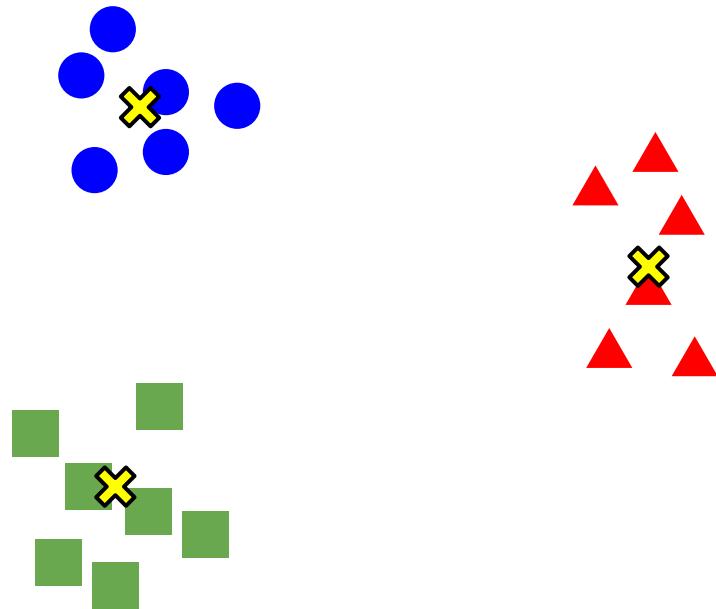
# k-means



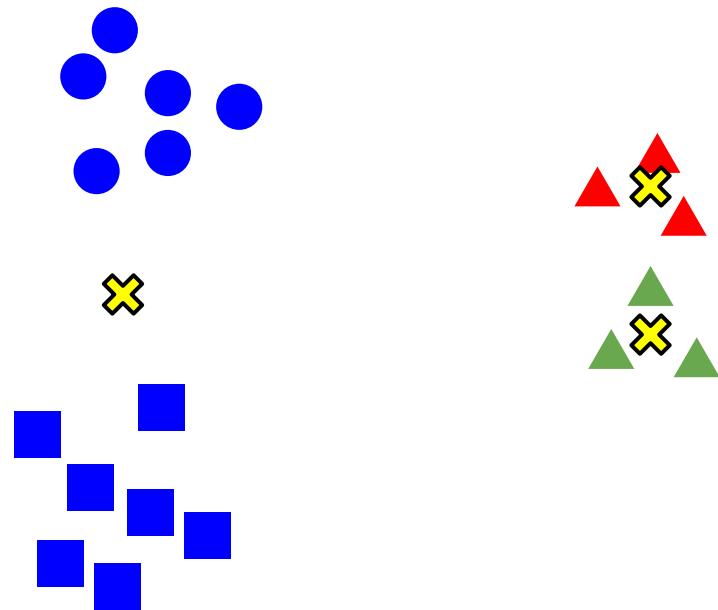
# k-means



# k-means

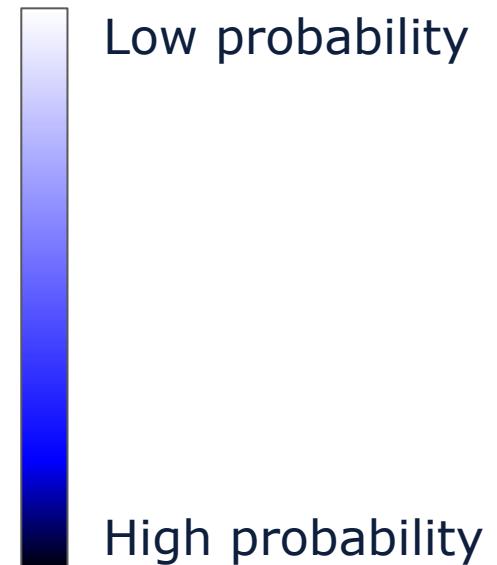
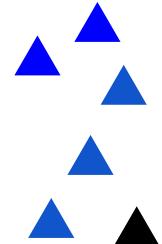
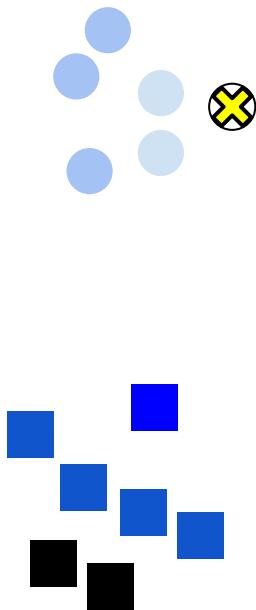


# Initial centers are important!



# k-means++

$d^2$  sampling

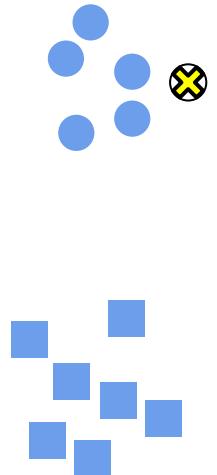


# A new family of clustering algorithms

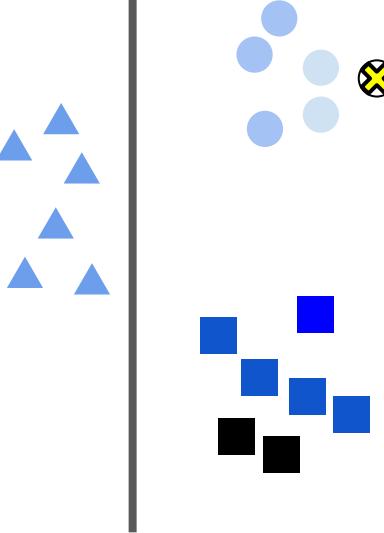
$d^\alpha$  sampling

$\beta$ -local search

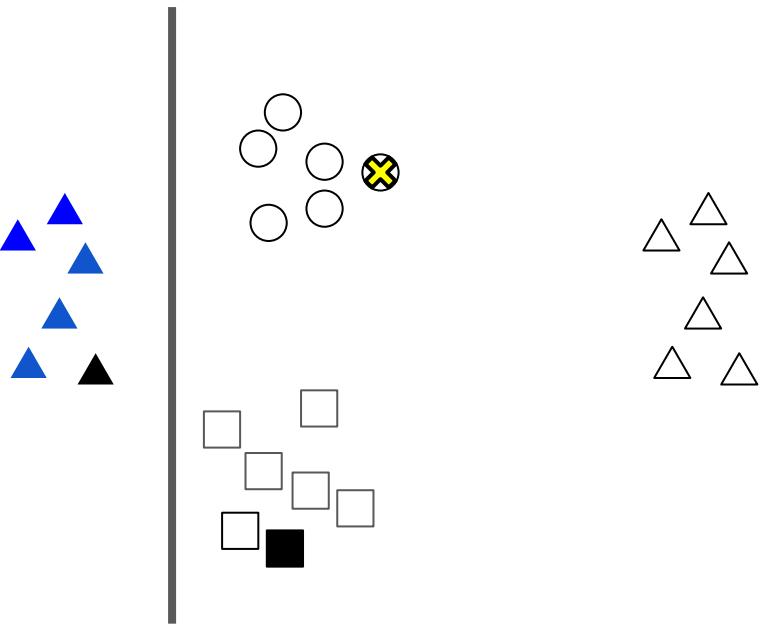
$\alpha = 0$



$\alpha = 2$



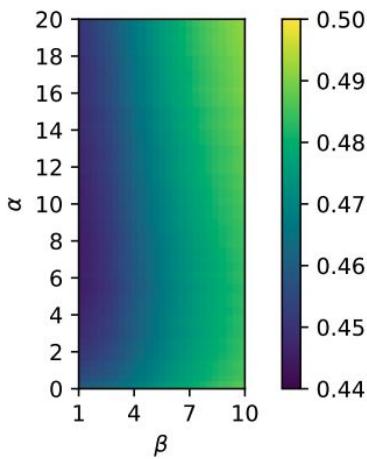
$\alpha = \infty$



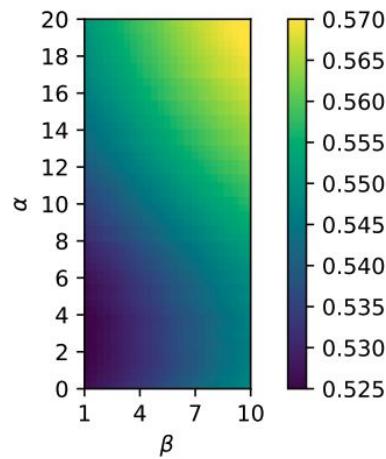
# A new family of clustering algorithms

$d^\alpha$  sampling

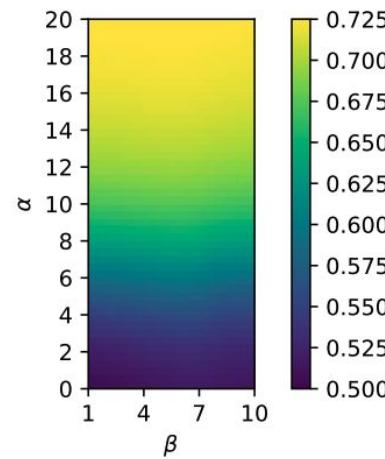
$\beta$ -local search



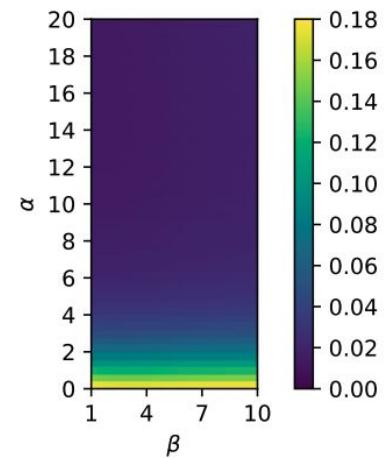
(a) MNIST



(b) CIFAR-10



(c) CNAE-9



(d) Gaussian Grid

Theorem: Given  $\tilde{O}\left(\frac{k \log n}{\epsilon^2}\right)$  sampled clustering instances,  
with high probability for all  $\alpha, \beta$ ,  
 $|\text{Avg performance over training set} - \text{expected performance}| < \epsilon$ .

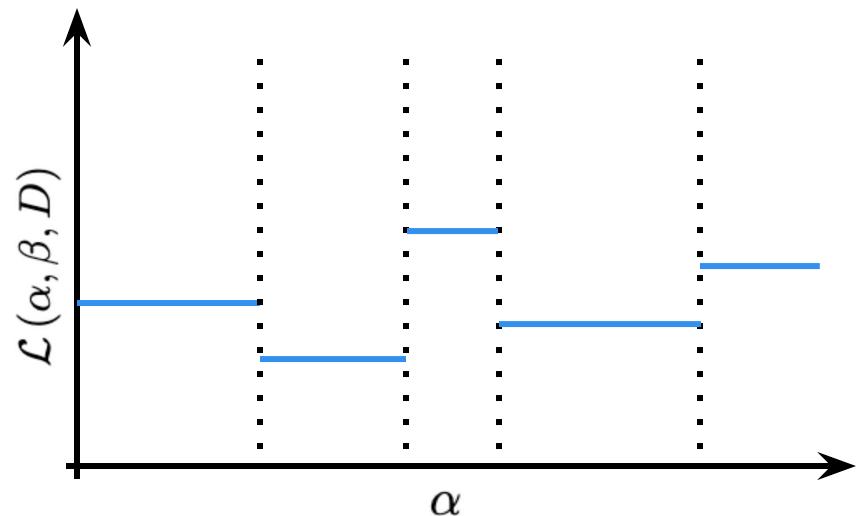
Key intuition:

A set of hyperparameters with high  
performance **on the training set**, also  
achieves high performance **in the future**

Theorem: Given  $\tilde{O}\left(\frac{k \log n}{\epsilon^2}\right)$  sampled clustering instances,  
 with high probability for all  $\alpha, \beta$ ,  
 $|\text{Avg performance over training set} - \text{expected performance}| < \epsilon$ .

Key insight:

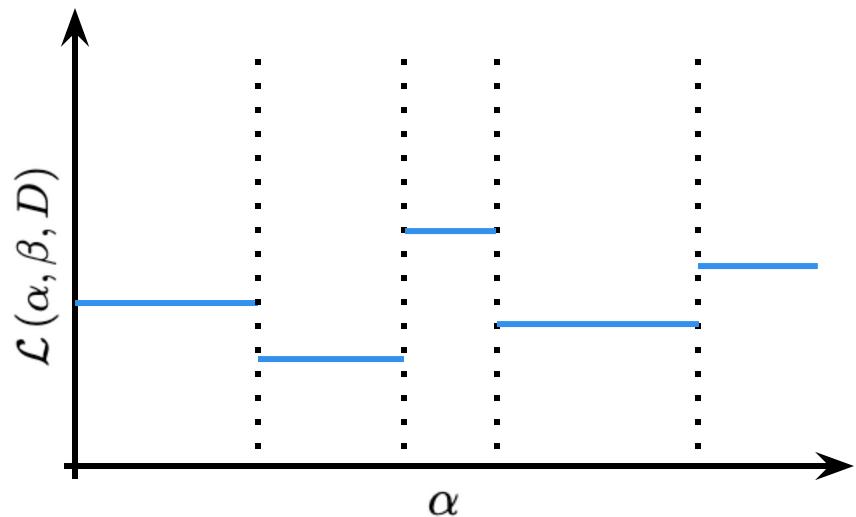
Bound the expected number  
 of **discontinuities** of  $\mathcal{L}(\alpha, \beta, D)$   
 as a function of  $\alpha, \beta$



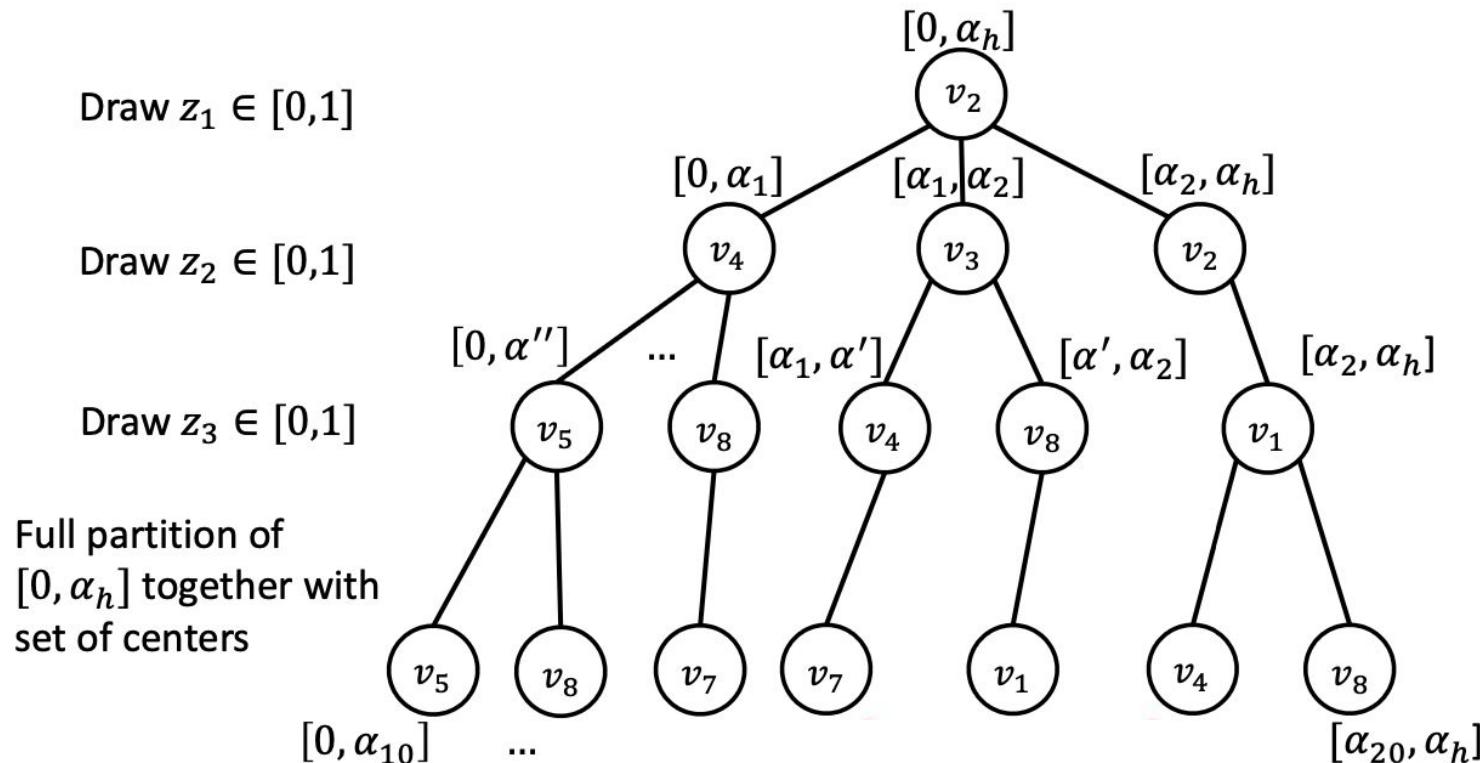
Theorem: Given  $\tilde{O}\left(\frac{k \log n}{\epsilon^2}\right)$  sampled clustering instances,  
with high probability for all  $\alpha, \beta$ ,  
 $|\text{Avg performance over training set} - \text{expected performance}| < \epsilon$ .

Efficient algorithm:

Solve for the **discontinuities**



# Execution Tree



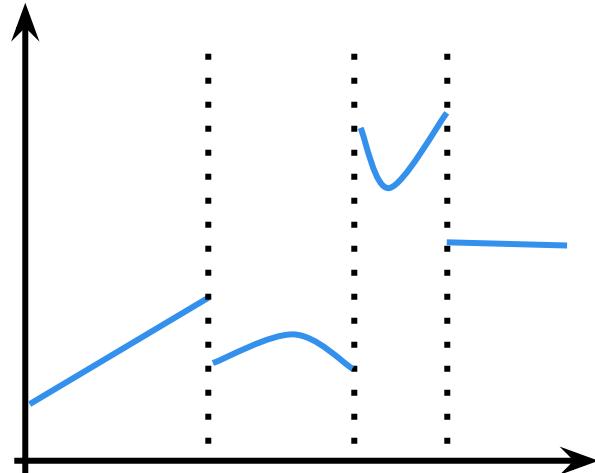
# Subsequent Work

## Generalizations

- Piecewise polynomial ([BDDKSV, 2019](#))
- Approximately piecewise polynomial ([BSV, 2020](#))

## Beyond clustering

- Integer Programming ([BDV, 2018](#))
- ElasticNet ([BKST, 2022](#))

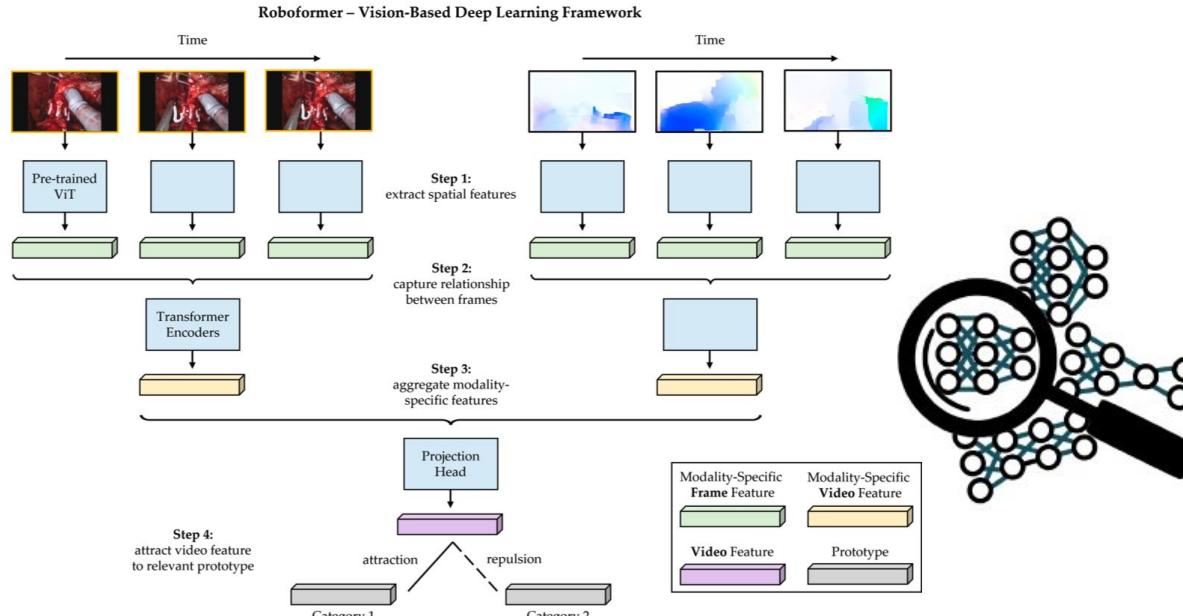


# Outline

- Introduction
- NAS: algorithm design
- HPO: theoretical results
- **Future directions**
  - Surgical AI
  - Climate sciences



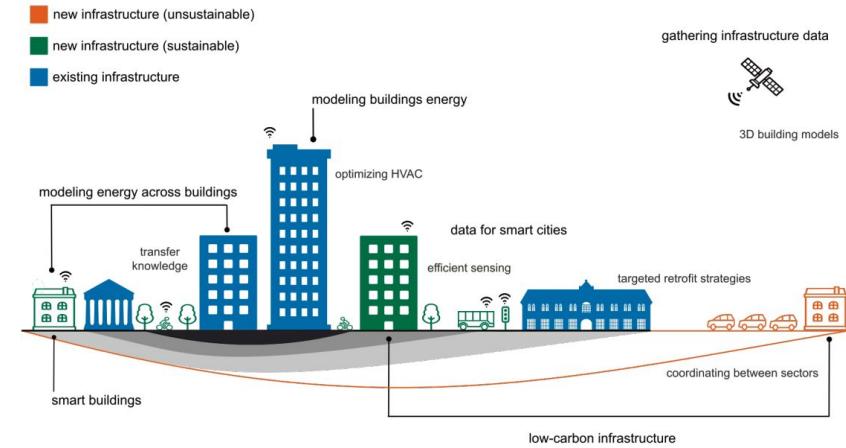
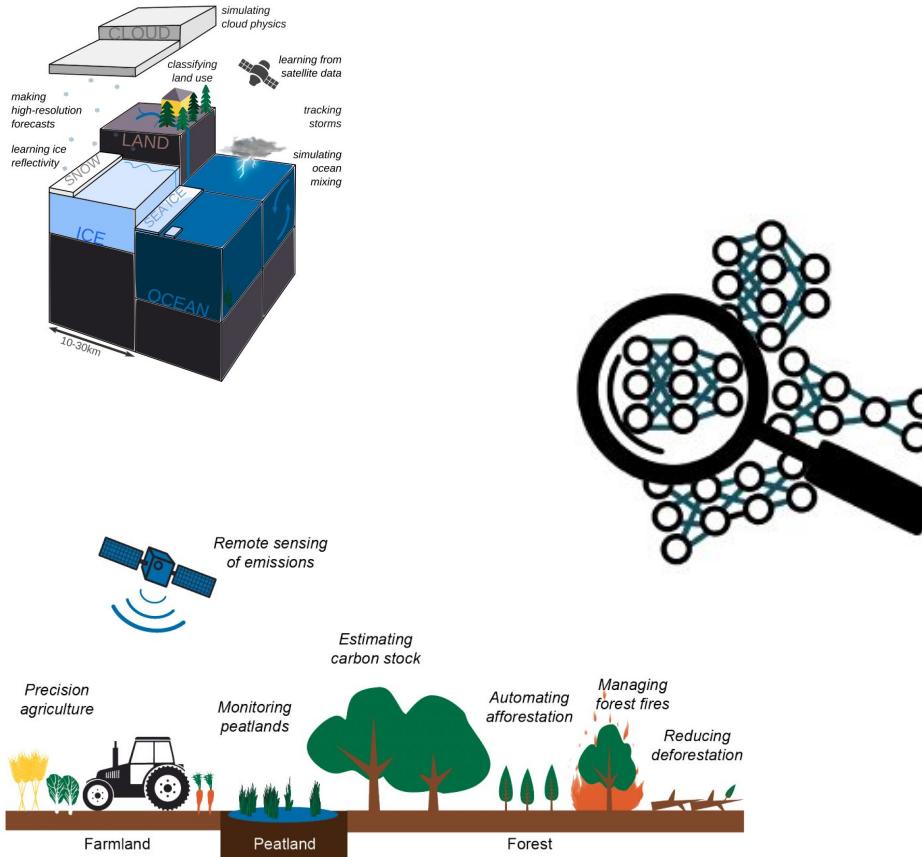
# Tackling Surgical AI with AutoML



[Kiyasseh et al. \(2022\)](#)

De-biasing surgical AI models with AutoML

# Tackling Climate Change with **AutoML**



[Rolnick et al. \(2018\)](#)  
[TRPNJSRTN W \(2022\)](#)

# Thanks!



Slides (with hyperlinks): <https://crwhite.ml/>