

# HOTNAS: Hierarchical Optimal Transport for Neural Architecture Search

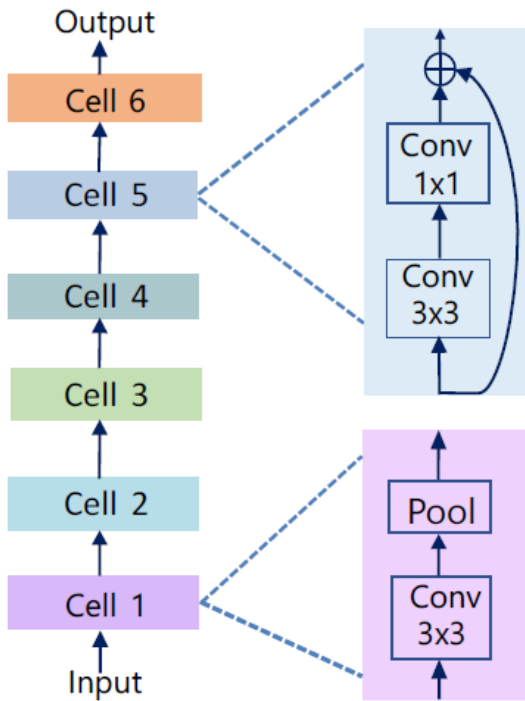
Poster: **WED-AM-359**

Jiechao Yang<sup>1,2</sup>   Yong Liu<sup>1,2</sup>   Hongteng Xu<sup>1,2</sup>

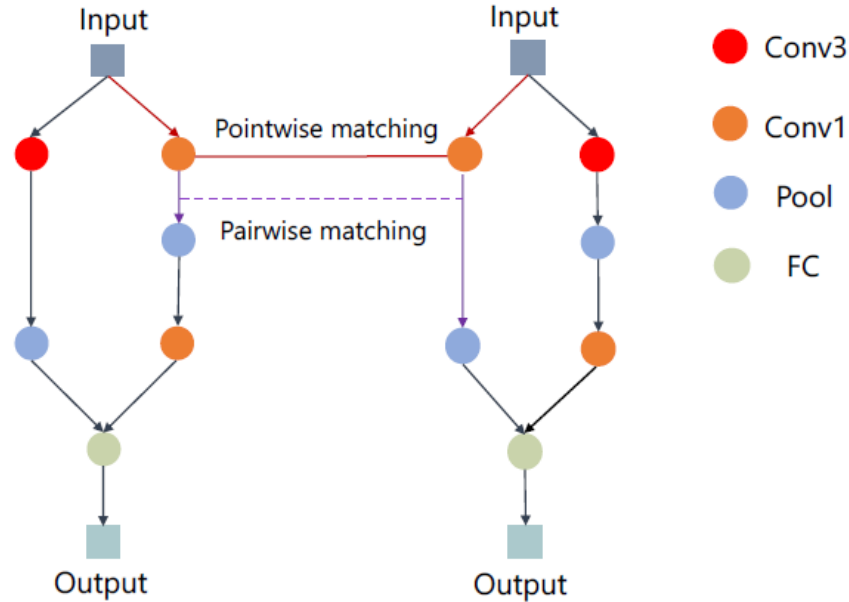
<sup>1</sup>Gaoling School of Artificial Intelligence, Renmin University of China, Beijing, China

<sup>2</sup>Beijing Key Laboratory of Big Data Management and Analysis Methods

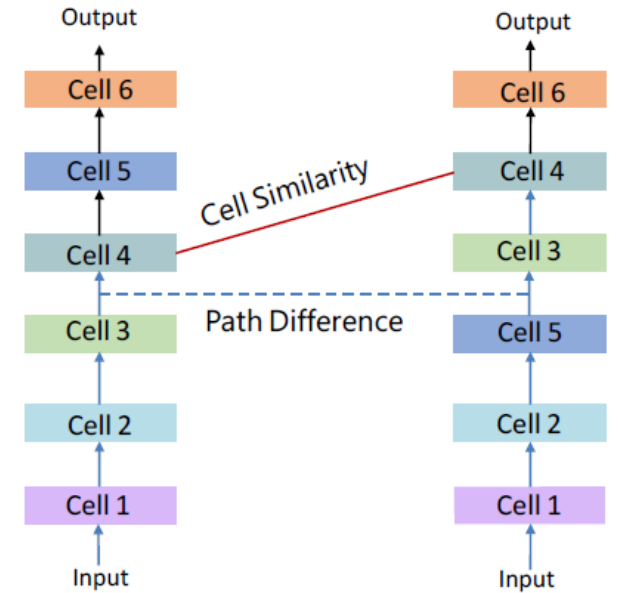
## HOTNAS: Hierarchical Optimal Transport for Neural Architecture Search



A modular cell-based network with hierarchical structure



(a) cell-level similarity



(b) network-level similarity

a **hierarchical optimal transportation metric** HOTNN , which jointly measures the similarity of cell internal architectures and the difference in macroarchitectures.

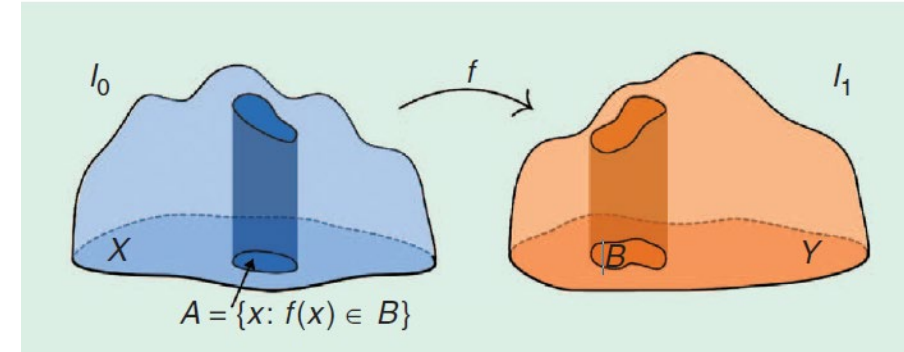
- The objective of NAS is to discover an optimal neural network architecture with the minimum validation loss

$$a^* = \operatorname{argmin}_{a \in \mathcal{A}} f(a),$$

- **Bayesian optimization (BO)** can **quickly** discover high-performing network architectures **with a limited number of samples**.

- The core of BO is **accurately measuring the similarity between different networks**.

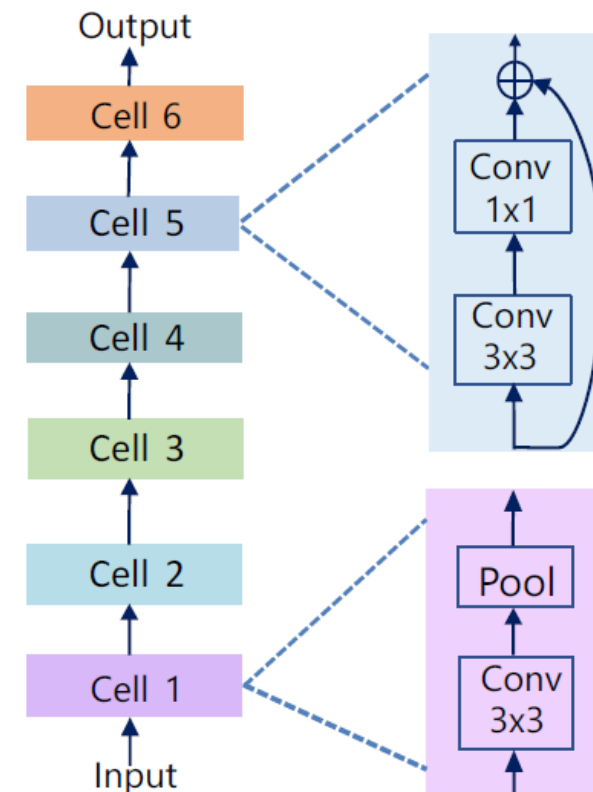
- Each network can be viewed as a **directed acyclic attributed graph**.
- **Optimal Transport (OT)** can naturally handle the **graph-like architecture**.
- NASBOT [Kandasamy et al. 2018] compute the minimum OT distance between networks as the similarity metric, but ignoring the similarity between cells.
- TW [Nguyen et al. 2021] are limited to searching for a single cell architecture and ignore the similarity of macro-architectures.



An example of OT [Kolouri et al. 2017]

## Network similarity metric: HOTTN

- measure the overall similarity between cell-based networks by leveraging its **hierarchical structure**
- organize the architecture into layers according to cells and learn the similarity within and between different layers.
- **cell-level similarity** computes the **OT distance between cells** in various networks by considering **the similarity of each node** and the **differences in the information flow costs between node pairs within each cell** in terms of **operational and structural information**.
- **Network-level similarity** calculates OT distance between networks by considering both **the cell-level similarity** and **the variation in the global position of each cell** within their respective networks



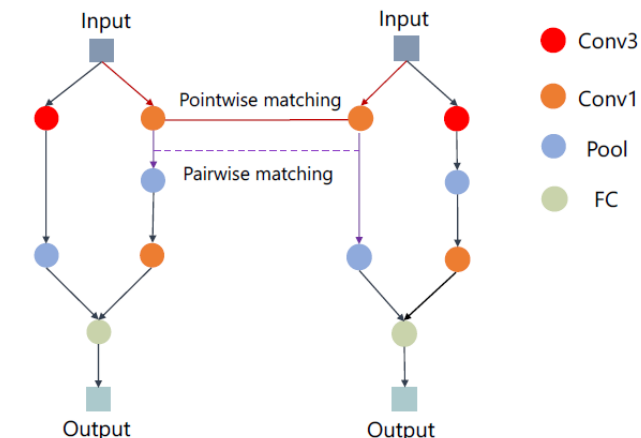
Example of a modular cell-based network

Two cell networks  $\mathcal{G}_{B_1} = (\mathcal{L}^{B_1}, \mathcal{E}^{B_1}, \ell_o^{B_1}, \ell_s^{B_1}), \quad | \quad \mathcal{G}_{B_2} = (\mathcal{L}^{B_2}, \mathcal{E}^{B_2}, \ell_o^{B_2}, \ell_s^{B_2})$

Probability measure  $\alpha = \sum_{i=1}^n p_i \delta(o_i^p, s_i^p) \quad \beta = \sum_{j=1}^m q_j \delta(o_j^q, s_j^q)$

Transport matrix  $\mathbf{T} \in U(\mathbf{p}, \mathbf{q}) = \{\mathbf{T} \in \mathbb{R}_+^{n \times m} \mid \mathbf{T}\mathbf{1}_m = \mathbf{p}, \mathbf{T}^T\mathbf{1}_n = \mathbf{q}\}$ ,

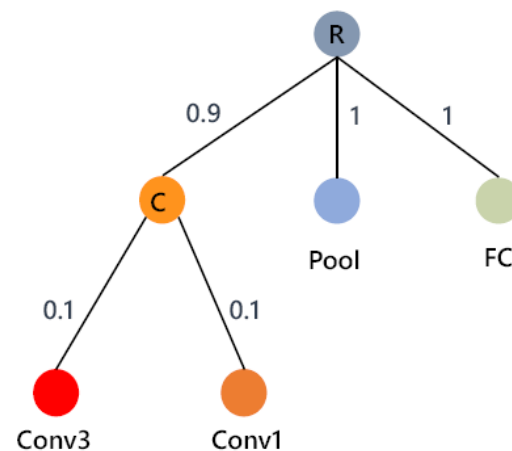
**Pointwise matching:** considering differences in the operational type and structural location information of each node.



$$\min_{\mathbf{T} \in U(\mathbf{p}, \mathbf{q})} \sum_{i=1}^n \sum_{j=1}^m \mathbf{T}_{ij} C_{ij}^{pq}$$

$$C_{ij}^{pq} = \varepsilon D(o_i^p, o_j^q) + (1 - \varepsilon) H(s_i^p, s_j^q), \quad 1 \leq i \leq n, 1 \leq j \leq m,$$

$$H(s_i^p, s_j^q) = \frac{1}{6} (s_i^p - s_j^q) \mathbf{1}_6,$$



$D(o_i^p, o_j^q)$

	Conv3	Conv1	Pool	FC
Conv3	0	0.1	2	2
Conv1	0.1	0	2	2
Pool	2	2	0	2
FC	2	2	2	0

■ **Pairwise matching** : learn the differences in the movement cost of various information flows between pairs of nodes within each cell network in terms of operational and structural information.

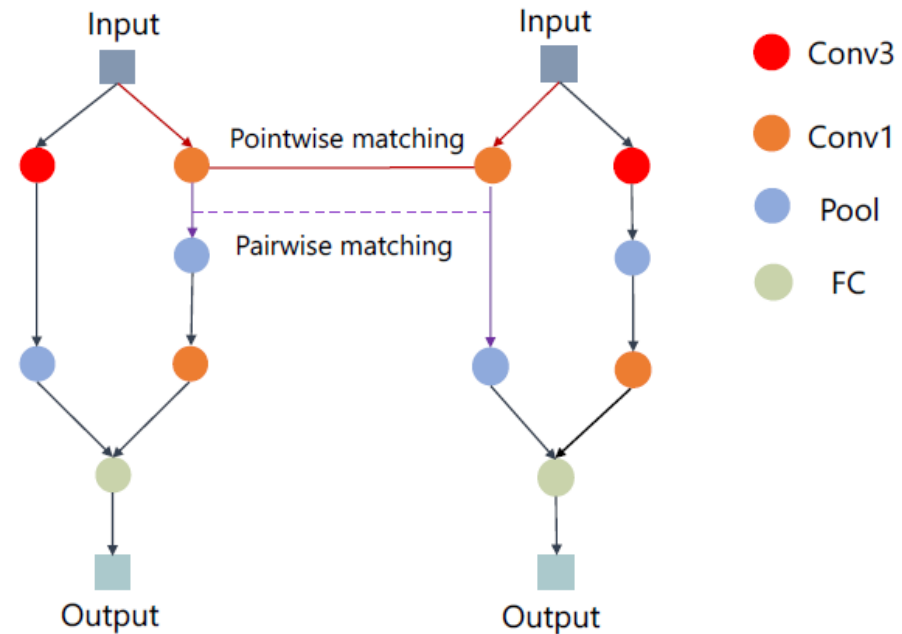
$$\min_{\mathbf{T} \in U(\mathbf{p}, \mathbf{q})} \sum_{i=1}^n \sum_{k=i+1}^n \sum_{j=1}^m \sum_{l=j+1}^m \mathbf{T}_{ij} \mathbf{T}_{kl} |C_{ik}^p - C_{jl}^q|.$$

$$C_{ik}^p = \varepsilon D(o_i^p, o_k^p) + (1 - \varepsilon) H(s_i^p, s_k^p), 1 \leq i < k \leq n,$$

$$C_{jl}^q = \varepsilon D(o_j^q, o_l^q) + (1 - \varepsilon) H(s_j^q, s_l^q), 1 \leq j < l \leq m.$$

■ **iFGW metric**

$$\begin{aligned} \text{iFGW}(\mathcal{G}_{B_1}, \mathcal{G}_{B_2}) = & \min_{\mathbf{T} \in U(\mathbf{p}, \mathbf{q})} \sum_{i=1}^n \sum_{k=i+1}^n \sum_{j=1}^m \sum_{l=j+1}^m \lambda \mathbf{T}_{ij} C_{ij}^p \\ & + (1 - \lambda) \mathbf{T}_{ij} \mathbf{T}_{kl} |C_{ik}^p - C_{jl}^q|, \end{aligned}$$



■ **Transport matching matrix**

$$\Gamma \in V(\mathbf{f}, \mathbf{g}) = \{\Gamma \in \mathbb{R}_+^{N \times M} \mid \Gamma \mathbf{1}_M = \mathbf{f}, \Gamma^T \mathbf{1}_N = \mathbf{g}\},$$

- **Cost matrix:** consider both the similarity between cells in two networks and the difference in the global position of each cell in their respective networks.

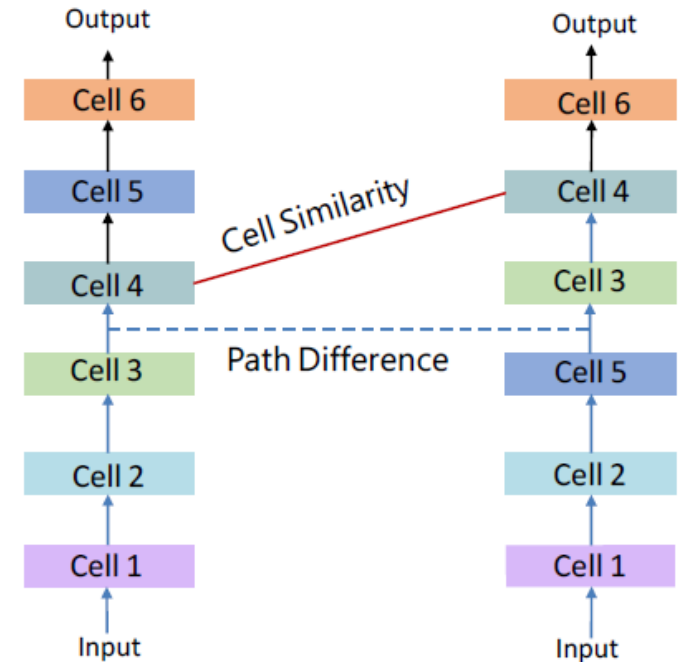
$$P(B_s^1, B_t^2) = \left| \frac{\delta^1(B_s^1)}{\delta^1} - \frac{\delta^2(B_t^2)}{\delta^2} \right|,$$

$$1 \leq s \leq N, 1 \leq t \leq M,$$

$$\mathbf{S}_{st}^{12} = (1 - \eta) \text{iFGW}(B_s^1, B_t^2) + \eta P(B_s^1, B_t^2)$$

■ **HOTNN metric**

$$\text{HOTNN}(\mathbf{a}^1, \mathbf{a}^2) = \min_{\Gamma \in V(\mathbf{f}, \mathbf{g})} \sum_{s=1}^N \sum_{t=1}^M \Gamma_{st} \mathbf{S}_{st}^{12}.$$



**Algorithm 1:** Hierarchical Optimal Transport for Neural Architecture Search

**Input:** Total number of iterations  $N$ , initial datapoints  $\mathcal{D}_0$ , search space  $\mathcal{A}$ , The maximum iterations  $T$

**Output:** The best architecture  $\mathbf{a}^*$

- 1 **for**  $t = 0$  **to**  $T - 1$  **do**
- 2     Compute HOTNN metric between different architectures on the current observation set  $\mathcal{D}_t = \mathcal{D}_0$ ;
- 3     Embed the HOTNN metric to the kernel function of GP;
- 4     Fit the GP on the current observation set  $\mathcal{D}_t$ ;
- 5     Construct the UCB acquisition function based on the predictive mean and variance (see Eq. (S7));
- 6     Generate a pool of candidate architectures  $\mathcal{P}_t$  by mutating the current best-performing architectures;
- 7     Select the next promising architectures  $\mathbf{a}_{\text{new}} = \arg \max_{\mathbf{a} \in \mathcal{P}_t} u_t(\mathbf{a})$ ;
- 8     Evaluate  $\mathbf{a}_{\text{new}}$  to obtain its validation loss  $y_{\text{new}}$ ;
- 9     Update the observation set  $\mathcal{D}_{t+1} = \mathcal{D}_t \cup \{\mathbf{a}_{\text{new}}, y_{\text{new}}\}$ ;
- 10 **end**
- 11 **return** the best-performing architecture  $\mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathcal{D}_T} f(\mathbf{a})$

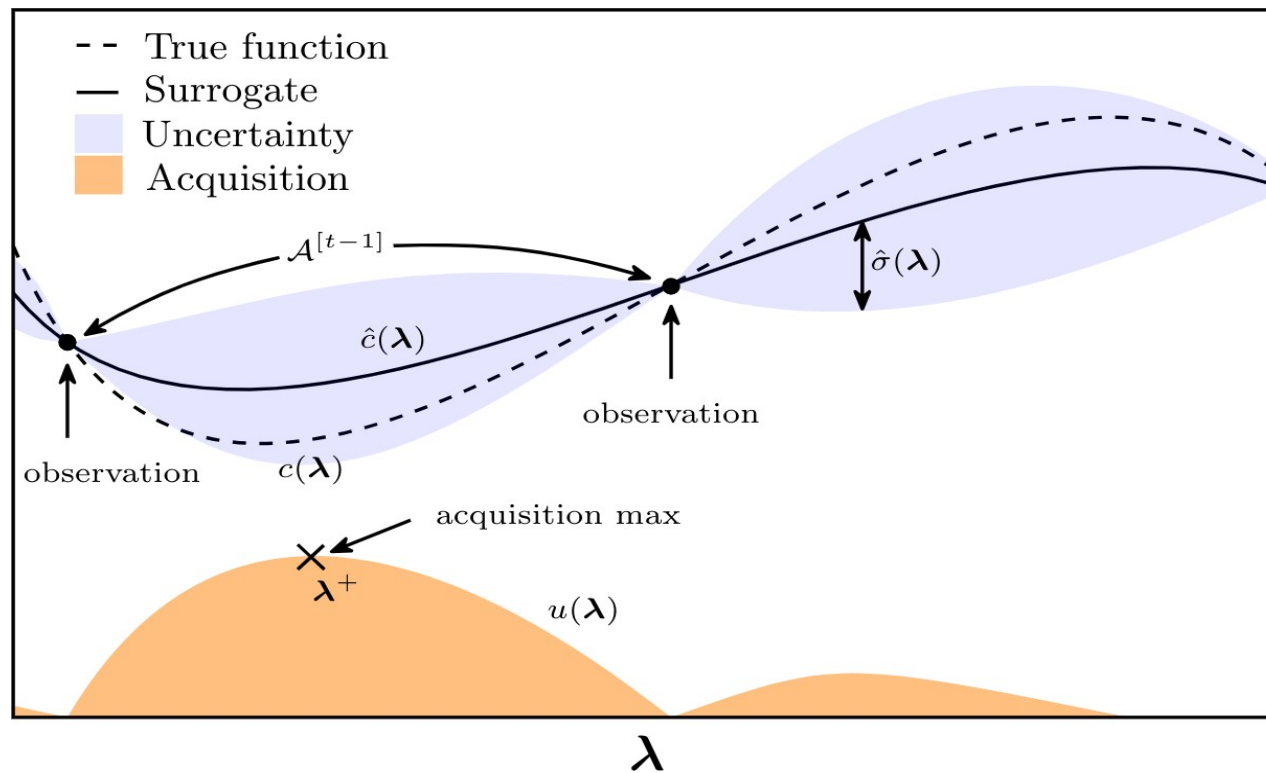


Illustration of BO



Table 1. Comparisons of the best-found valid loss and test loss on the TransNAS-Bench-101 benchmark and the DARTS benchmark.

Search Space	Tasks	Loss	Random Search	Evolutionary Search	BO-edit	NASBOT	HOTNAS
TransNAS-Bench-101	Autoencoding	Valid Error	28.09±0.18	28.19±0.24	28.55±0.28	29.25±0.25	<b>25.80±0.04</b>
		Test Error	26.65±0.18	26.74±0.24	27.14±0.28	27.82±0.25	<b>24.36±0.01</b>
	Object Classification	Valid Error	53.21±0.02	52.96±0.02	53.42±0.03	53.28±0.02	<b>52.69±0.00</b>
		Test Error	46.49±0.02	46.25±0.04	46.67±0.03	46.30±0.03	<b>45.90±0.02</b>
	Scene Classification	Valid Error	43.85±0.03	43.43±0.03	43.58±0.03	43.78±0.04	<b>43.19±0.01</b>
		Test Error	35.43±0.02	35.25±0.03	35.29±0.02	35.25±0.03	<b>35.00±0.01</b>
	Jigsaw	Valid Error	3.58±0.03	3.25±0.01	3.31±0.00	3.27±0.01	<b>3.17±0.01</b>
		Test Error	3.79±0.04	3.35±0.01	3.34±0.01	3.38±0.02	<b>3.29±0.01</b>
	Surface Normal	Valid Error	39.20±0.04	37.55±0.16	37.42±0.14	38.80±0.15	<b>36.65±0.20</b>
		Test Error	36.27±0.04	34.72±0.15	34.69±0.14	35.99±0.15	<b>33.91±0.19</b>
	Room Layout	Valid Error	59.98±0.04	59.83±0.03	59.95±0.06	60.19±0.05	<b>58.92±0.05</b>
		Test Error	53.94±0.06	55.54±0.10	54.02±0.03	54.72±0.09	<b>53.80±0.04</b>
	Semantic Segmentation	Valid Error	71.59±0.04	71.39±0.05	71.01±0.05	70.89±0.04	<b>70.51±0.01</b>
		Test Error	68.97±0.01	68.11±0.05	68.45±0.05	68.30±0.04	<b>67.98±0.03</b>
DARTS	CIFAR-10	Valid Error	5.90±0.07	5.50±0.09	5.42±0.14	5.73±0.07	<b>5.37±0.01</b>
		Test Error	3.28±0.09	2.87±0.04	2.72±0.07	2.93±0.12	<b>2.43±0.04</b>
	CIFAR-100	Test Error	21.47±0.08	19.75±0.13	20.62±0.12	19.95±0.17	<b>18.46±0.09</b>

# Experiments

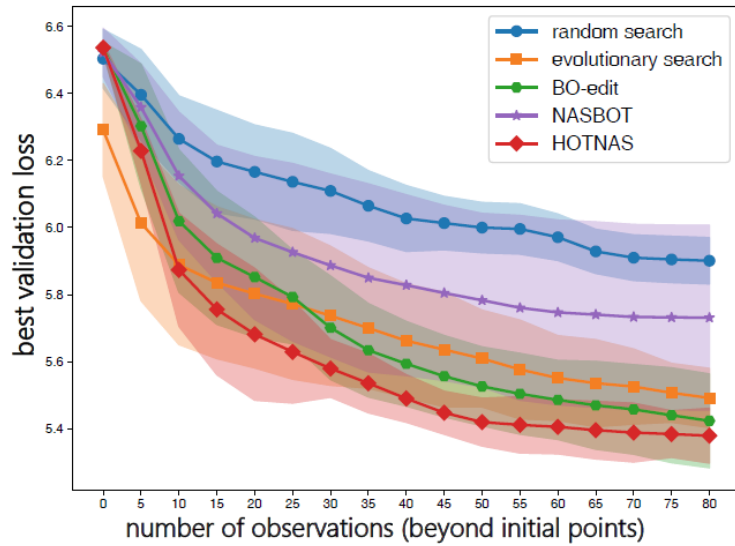
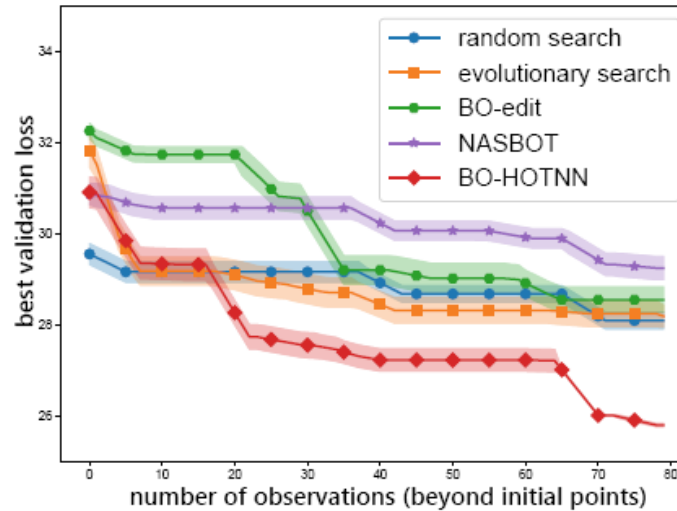
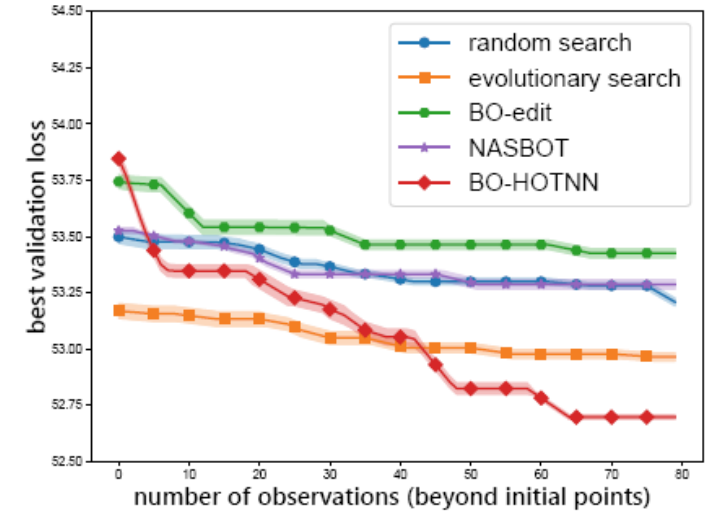


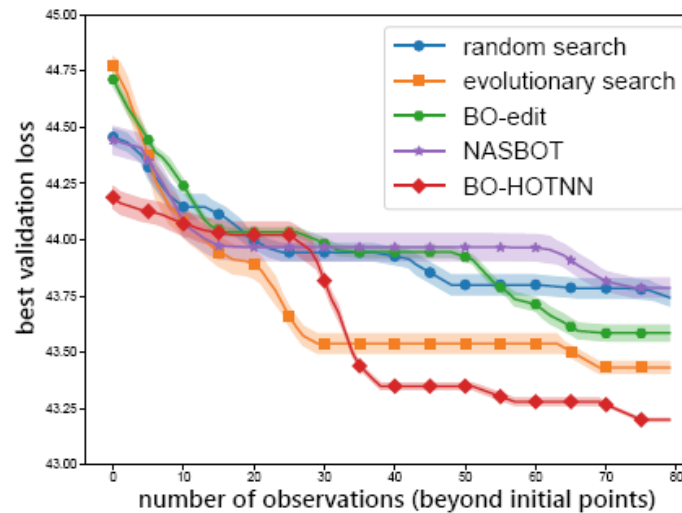
Figure 3. The best-found validation loss over the number of iterations (beyond initial points) of various NAS methods on the DARTS benchmark.



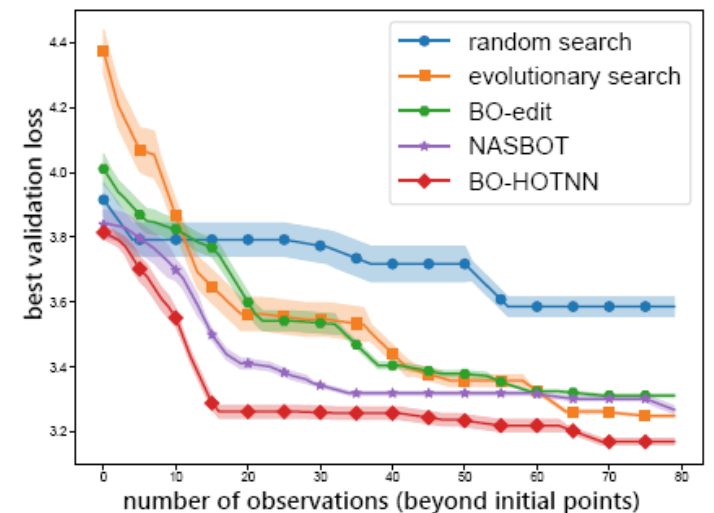
(a) Autoencoding



(b) Object classification



(c) Scene classification



(d) Jigsaw



# Thanks!