

IT-DSE: Invariance Risk Minimized Transfer Microarchitecture Design Space Exploration

Ziyang Yu¹, Chen Bai¹, Shoubo Hu², Ran Chen², Taohai He³,
Mingxuan Yuan², **Bei Yu**¹, Martin Wong¹

¹ The Chinese University of Hong Kong

² Huawei Noah's Ark Lab

³ HiSilicon

Oct. 31, 2023



- 1 Background
- 2 Preliminaries
- 3 IT-DSE Algorithm
- 4 Experiments

Background

Microarchitecture Design Space exploration:

- Microarchitecture determines the detailed structures of a microprocessor.
- Microarchitecture design space exploration involves finding different configurations with desired performance, power and area (PPA).

Two major challenges:

- Complexity of design spaces growing rapidly: number of potential configurations, inter-feature interactions, mixed-type parameters.
- Getting desired evaluation metric reports for one microarchitecture configuration through VLSI integration flow is timing consuming.

Previous solution

In industry:

- Manually configure design parameters by computer architects.
- **Limitation:** Requires both sufficient domain knowledge and a large amount of human labor for computer architects.

In academia:

- Machine learning-based process simulation models or surrogate models.
- **Limitations:** Effectiveness depends on quality of training data and model granularity.

Two facts:

- Different generations of microprocessors are developed based on a common baseline microarchitecture. Minor improvements are applied across consecutive microarchitecture generations.
- Some general relations exist for design features due to the similar microarchitecture design paradigm.

ANN-TL¹: Leveraging historical configurations for new tasks with cross-domain mixup and Artificial Neural Network (ANN) feature.

Limitation: ignores intricate feature interactions, lacks interpretability, presumes identical design spaces for source and target tasks.

¹Jianwang Zhai, Yici Cai, and Bei Yu (2023). “Microarchitecture Power Modeling via Artificial Neural Network and Transfer Learning”. In: *2023 28th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 1–6.

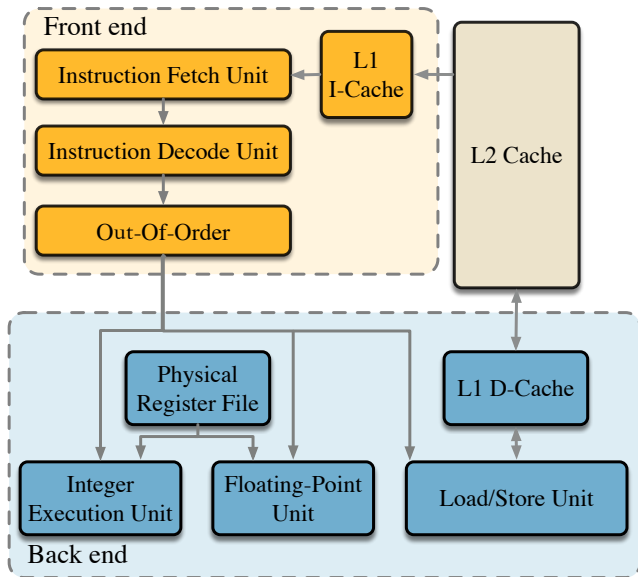
The characteristics of certain design features have **similar or even identical effects** on PPA values in different generations of a microarchitecture design.

Extract the relationship that is **approximately invariant** among different historical source design tasks, even equipped with different design spaces.

Transfer domain knowledge from different source tasks to the target design space exploration task.

Preliminaries

Typical Microprocessor Components



Visualization of typical microarchitecture pipeline.

Prediction function $f' : \mathcal{X} \rightarrow \mathcal{Y} \Rightarrow \phi_{\mathbf{u}}(\cdot)$: feature extractor; $h_{\mathbf{w}}(\cdot)$: regressor.
 \mathbf{u}, \mathbf{w} : parameter.

IRM:

$$\begin{aligned} \min_{\mathbf{u}, \mathbf{w}} \quad & \sum_{s=1}^S \mathcal{R}^s(\mathbf{u}, \mathbf{w}), \\ \text{s.t.} \quad & \mathbf{w} \in \underset{\mathbf{w}^s}{\operatorname{argmin}} \mathcal{R}^s(\mathbf{u}, \mathbf{w}^s), \end{aligned} \tag{1}$$

$\mathcal{R}^s(\mathbf{u}, \mathbf{w}) = \mathbb{E}_{X^s, Y^s}[\mathcal{L}(f'(X^s), Y^s)]$: risk under s -th source task.
 $\mathcal{L}(\cdot)$: the loss function.

Simplified IRMv1:

$$\min_{\mathbf{u}, \mathbf{w}} \sum_{s=1}^S \mathcal{R}^s(\mathbf{u}, \mathbf{w}) + \lambda \|\nabla_{\mathbf{w}} \mathcal{R}^s(\mathbf{u}, \mathbf{w})\|^2, \tag{2}$$

²Martin Arjovsky et al. (2019). "Invariant risk minimization". In: *arXiv preprint arXiv:1907.02893*. [10/30](https://arxiv.org/abs/1907.02893)

Source Task

Previously explored microarchitecture designs tasks.

Composed of the explored sample dataset \mathcal{D}^s containing n_s parameter vectors

$X^s = [\mathbf{x}_1^s, \dots, \mathbf{x}_{n_s}^s]^\top \in \mathbb{R}^{n_s \times d_s}$ and the evaluated PPA vectors $Y^s = [\mathbf{y}_1^s, \dots, \mathbf{y}_{n_s}^s] \in \mathbb{R}^{n_s \times 3}$.

Target Task

The new microarchitecture design space exploration task with only the set of legal parameter configuration vectors $X^t = [\mathbf{x}_1^t, \dots, \mathbf{x}_{n_t}^t]^\top \in \mathbb{R}^{n_t \times d_t}$ available.

Pareto Optimality

For a multi-objective minimization problem with M objectives, \mathbf{x}_1 is deemed to dominate solution \mathbf{x}_2 ($\mathbf{x}_1 \succeq \mathbf{x}_2$) if

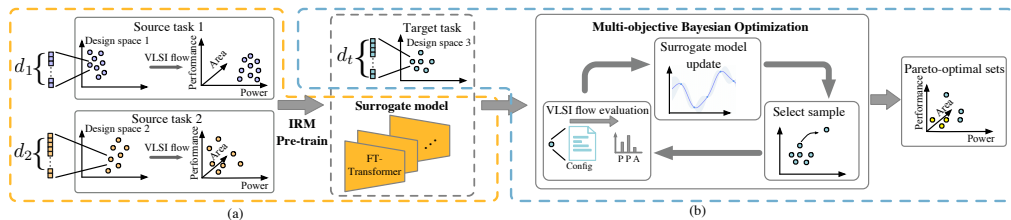
$$\begin{aligned} f_m(\mathbf{x}_1) &\leq f_m(\mathbf{x}_2), \forall m \in \{1, \dots, M\}, \\ f_k(\mathbf{x}_1) &< f_k(\mathbf{x}_2), \exists k \in \{1, \dots, M\}. \end{aligned} \tag{3}$$

Pareto-optimal set: The collection of solutions that remain non-dominated by others.

Microarchitecture Transferring Design Space Exploration

Given S source tasks with explored datasets $\{\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^S\}$ and the target sample \mathcal{D}^t , the objective is to utilize the information from historical source tasks and improve the efficiency of finding a series of microprocessor design configurations X in the target task that forms the Pareto optimality among the associated subset of $Y \subset \mathcal{Y}$, so that $X = \{x | x \succeq x', \forall x' \in X^t\}$, $Y = \{f(x) | X \in X^t\}$.

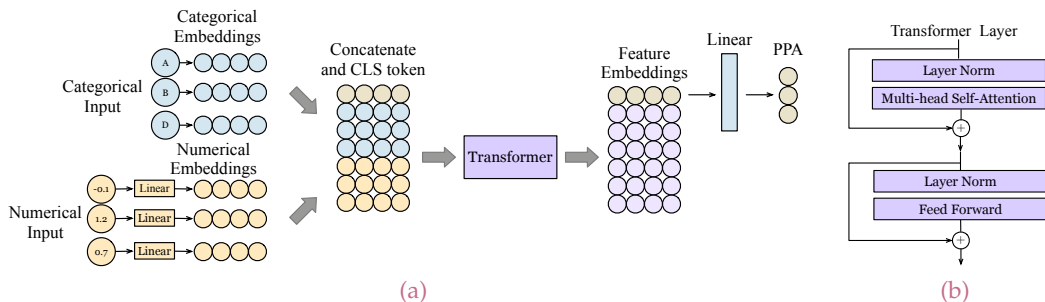
IT-DSE Algorithm



Workflow of proposed IT-DSE.

- Accommodate varying design spaces: FT-Transformer ensemble
- Knowledge fusion and warm start new task: pre-training with invariant risk minimization paradigm

Customized Surrogate Model with Transformer



(a) FT-Transformer architecture. Firstly, Feature Tokenizer transforms input parameters into embeddings. The embedding is then processed by the Transformer module. (b) One Transformer layer.

$$\mathbf{t}_i^{num} = x_i^{num} \cdot \mathbf{w}_i^{num} + b_i^{num}, \quad (4)$$

$$\mathbf{t}_j^{cat} = \mathbf{e}_j^\top \mathbf{w}_j^{cat} + b_j^{cat}, \quad (5)$$

$$\mathbf{T} = \text{stack}([\text{CLS}], \mathbf{t}_1^{num}, \dots, \mathbf{t}_{d_n}^{num}, \mathbf{t}_1^{cat}, \dots, \mathbf{t}_{d_c}^{cat}). \quad (6)$$

Comprehensive union of extracted features:

Source task 1: $\mathbf{x}^1 = [x_1, x_2]$, source task 2: $\mathbf{x}^2 = [x_2, x_3, x_4]$,

Embedding weight: $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3, \mathbf{w}_4]^\top \in \mathbb{R}^{4 \times d}$,

Embedding bias: $\mathbf{B} = [b_1, b_2, b_3, b_4]^\top \in \mathbb{R}^{4 \times d}$,

Transformed embeddings: $\mathbf{T}^1 = \text{stack}([\text{CLS}], \mathbf{t}_1, \mathbf{t}_2)$, $\mathbf{T}^2 = \text{stack}([\text{CLS}], \mathbf{t}_2, \mathbf{t}_3, \mathbf{t}_4)$.

Focusing more on task-invariant features.

For s -th source task, original data: $\mathcal{D}^s = \{X^s, Y^s\} = \{\mathbf{x}_i^s, \mathbf{y}_i^s\}_{i=1}^{n_s}$,

transferred data: $\mathcal{T}^s = \{\phi_u(\mathbf{x}_i^s), \mathbf{y}_i^s\}_{i=1}^{n_s}$,

collection of transformed datasets of all S source tasks: $\mathcal{T}^c = \cup_{s=1}^S \mathcal{T}^s$.

Objective for Bayesian IRM:

$$\begin{aligned} \max_{\mathbf{u}} \sum_s \mathbb{E}_{g_u(w)} [\ln p(\mathcal{T}^s | w, \mathbf{u})] + \lambda \left(\mathbb{E}_{g_u(w)} [\ln p(\mathcal{T}^s | w, \mathbf{u})] \right. \\ \left. - \mathbb{E}_{g_u^s(w^s)} [\ln p(\mathcal{T}^s | w^s, \mathbf{u})] \right). \end{aligned} \quad (7)$$

$g_u(w) \approx p(w | \mathcal{T}^c)$, $g_u^s(w^s) \approx p(w^s | \mathcal{T}^s)$.

First term: Encouraging \mathbf{u} to retain information for data distribution fitting.

Second term: Requires transformed data distribution to be stable among different tasks with feature extractor $\phi_u(\cdot)$.

The actual posterior $p(\mathbf{w}|\mathcal{T}^c)$ and $p(\mathbf{w}^s|\mathcal{T}^s)$ are hard to estimate in large models. Variational inference: approximate the posterior distributions by maximizing the evidence lower bound (ELBO).

$$g_u^s(\mathbf{w}^s) = \operatorname{argmax}_{g' \in \mathcal{G}} \mathbb{E}_{g'}[\ln p(\mathcal{T}^s|\mathbf{w}, \mathbf{u}) - D_{\text{KL}}(g' || p_0(\mathbf{w}))], \quad (8)$$

$$g_u(\mathbf{w}) = \operatorname{argmax}_{g'} \sum_{s=1}^S \mathbb{E}_{g' \in \mathcal{G}}[\ln p(\mathcal{T}^s|\mathbf{w}, \mathbf{u}) - D_{\text{KL}}(g' || p_0(\mathbf{w}))], \quad (9)$$

Assume $g_u(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $g_u^s(\mathbf{w}^s)$ approaches $g_u(\mathbf{w})$ when more task-invariant features are extracted:

$$g_u^s(\mathbf{w}^s) = \mathcal{N}(\boldsymbol{\mu} - \nabla_{\boldsymbol{\mu}} \mathbb{E}_{g_u(\mathbf{w})} \ln p(\mathcal{T}^s|\mathbf{w}, \mathbf{u}), \boldsymbol{\Sigma}), \quad (10)$$

Use **Monte Carlo samples** to estimate penalty term in Equation (7):

$$\mathbf{w} = \boldsymbol{\mu} + \epsilon \boldsymbol{\Sigma}, \quad \mathbf{w}^s = \boldsymbol{\mu}^s + \epsilon \boldsymbol{\Sigma}^s. \quad (11)$$

Algorithm Feature extraction with Bayesian IRM

Require: Feature extractor ϕ_u , regressor h_w , collection of data from S source task $\{\mathcal{D}^s\}_{s=1}^S$

Ensure: the learned task-invariant feature extractor ϕ_u , regressor h_w .

- 1: Initialization: prior $\phi_0 \leftarrow \mathcal{N}(\mathbf{0}, \sigma \mathbf{I})$.
 - 2: **repeat**
 - 3: **for** s in $1, \dots, S$ **do**
 - 4: Sample a batch of data $(X_{batch}^s, Y_{batch}^s)$ from \mathcal{D}^s ;
 - 5: Transform data $\mathcal{T}_{batch}^s \leftarrow (\phi_u(X_{batch}^s), Y_{batch}^s)$;
 - 6: **end for**
 - 7: Update $g_u(w)$; ▷ Equation (9)
 - 8: **for** s in $1, \dots, S$ **do**
 - 9: Update $g_u^s(w^s)$; ▷ Equation (10)
 - 10: **end for**
 - 11: Sample $g_u(w)$ and $g_u^s(w^s)$; ▷ Equation (11)
 - 12: Update u to maximize objective; ▷ Equation (7)
 - 13: **until** Training finished.
 - 14: Return ϕ_u, h_w ;
-

Target Task DSE in Feature Space: Deep Ensemble Surrogate Model



Given n evaluated samples $\mathcal{D}^t = \{X^t, f'(X^t)\}$ from target task,
 $f'(X^t) = \{(f'(\mathbf{x}_1^t))^\top, f'(\mathbf{x}_2^t)^\top, \dots, f'(\mathbf{x}_n^t)^\top\}^\top \in \mathbb{R}^{n \times 3}$,

For new point \mathbf{x}_*^t , posterior distribution $p(\mathbf{y}_*^t | \mathbf{x}_*^t)$ is approximated as a **Gaussian distributions**:

$$p(\mathbf{y}_*^t | \mathbf{x}_*^t) = \mathcal{N} \left(\boldsymbol{\mu}_{\mathbf{y}_*^t | f'(X^t)}, \boldsymbol{\sigma}_{\mathbf{y}_*^t | f'(X^t)}^2 \right), \quad (12)$$

The posterior mean and variance:

$$\boldsymbol{\mu}_{\mathbf{y}_*^t | f'(X^t)} = \frac{1}{M} \sum_{m=1}^M \boldsymbol{\mu}_{\theta_m}(\mathbf{x}_*^t), \quad (13)$$

$$\begin{aligned} \boldsymbol{\sigma}_{\mathbf{y}_*^t | f'(X^t)}^2 &= \frac{1}{M} \sum_{m=1}^M (\boldsymbol{\mu}_{\mathbf{y}_*^t | f'(X^t)} - \boldsymbol{\mu}_{\theta_m}(\mathbf{x}_*^t))^2 \\ &+ \frac{1}{M} \sum_{m=1}^M \boldsymbol{\sigma}_{\theta_m}^2(\mathbf{x}_*^t). \end{aligned} \quad (14)$$

Target Task DSE in Feature Space: Pareto EHVI Acquisition Function³

Hypervolume given an approximate Pareto frontier $\mathcal{P}(\mathbf{y})$:

$$HV_{\mathbf{y}_{ref}}(\mathcal{P}(\mathbf{y})) = \lambda_M(\cup_{\mathbf{y} \in \mathcal{P}(\mathbf{y})} [\mathbf{y}, \mathbf{y}_{ref}]), \quad (15)$$

$[\mathbf{y}, \mathbf{y}_{ref}]$: hyperrectangle bounded by \mathbf{y} and reference point \mathbf{y}_{ref} .

Hypervolume improvement (HVI) for a new given point \mathbf{y}_* given $\mathcal{P}(\mathbf{y})$:

$$HVI(\mathbf{y}_* | \mathcal{P}(\mathbf{y}), \mathbf{y}_{ref}) = V_{\mathbf{y}_{ref}}(\mathcal{P}(\mathbf{y} \cup \mathbf{y}_*)) - V_{\mathbf{y}_{ref}}(\mathcal{P}(\mathbf{y})). \quad (16)$$

EHVI Monte Carlo integration approximation with K sampling times:

$$\begin{aligned} EHVI(\hat{\mathbf{x}}) &= \mathbb{E}_{P(f'(\hat{\mathbf{x}}_* | \mathcal{D}))} (HVI(f'(\hat{\mathbf{x}}) | \mathcal{P}(\mathbf{y}), \mathbf{y}_{ref})) \\ &\approx \frac{1}{n_k} \sum_{k=1}^K HVI(\tilde{f}'_k(\hat{\mathbf{x}}) | \mathcal{P}(\mathbf{y}), \mathbf{y}_{ref}). \end{aligned} \quad (17)$$

\tilde{f}'_k : k -th function evaluation from distribution $p(f'(\mathbf{x} | \mathcal{D}))$.

³Samuel Daulton, Maximilian Balandat, and Eytan Bakshy (2020). “Differentiable expected hypervolume improvement for parallel multi-objective Bayesian optimization”. In: *Advances in Neural Information Processing Systems* 33, pp. 9851–9864.

Experiments

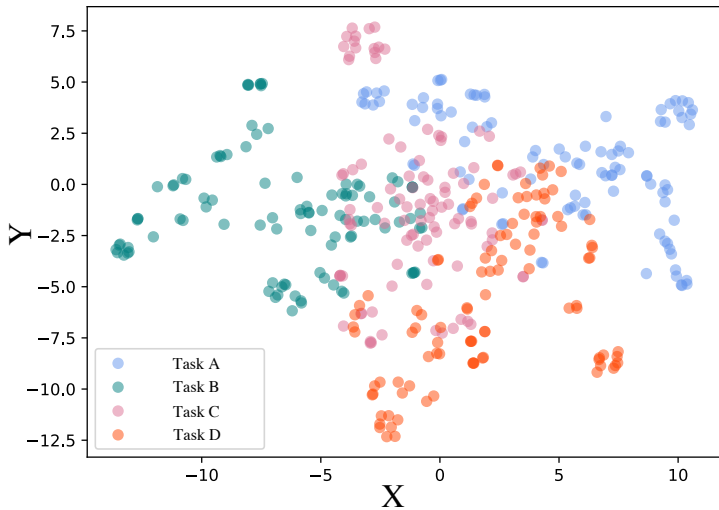
- Experiments are performed on an in-house 64-bit high-performance commercial microprocessor.
- VLSI flow is used to evaluate the PPA values of each microarchitecture configuration.
- Register-transfer-level (RTL) design: generated from Chisel⁴.
- Performance and power: obtained from the benchmark simulation after the physical implementation.
- Area value: reported from the physical implementation tool.
- Due to the commercial confidentiality of microarchitecture, we normalize the original PPA values.

⁴Jonathan Bachrach et al. (2012). “Chisel: Constructing Hardware in a Scala Embedded Language”. In: *ACM/IEEE Design Automation Conference (DAC)*, pp. 1216–1225.

Table: Statistics of our microarchitecture design space

Module	# Linear	#Pow	# Categorical	# Combinations
IFU	8	4	0	$\sim 4 \times 10^8$
OoO	11	0	0	$\sim 5 \times 10^9$
IEX	12	0	3	$\sim 1 \times 10^9$
FSU	5	0	0	$\sim 2 \times 10^3$
LSU	6	3	0	$\sim 1 \times 10^7$
L2C	1	2	0	$\sim 8 \times 10^2$
Overall	43	9	3	$\sim 3 \times 10^{40}$

- 55 parameters with more than 3×10^{40} combinations.
- Task A, task B, and Task C share the same design space, for task D, we can only tune 53 parameters out of 55 parameters listed.
- Task A, task B, task C, and Task D contain 1237, 377, 1835 and 3453 evaluated samples, respectively.



Data distribution for 4 tasks.

Hypervolume:

$$HV_{\mathbf{y}_{ref}}(\mathcal{P}(\mathbf{y})) = \lambda_M(\cup_{\mathbf{y} \in \mathcal{P}(\mathbf{y})} [\mathbf{y}, \mathbf{y}_{ref}]), \quad (18)$$

Average distance to reference set (ADRS):

$$ADRS(\mathcal{P}^*, \mathcal{P}) = \frac{1}{|\mathcal{P}^*|} \sum_{\alpha \in \mathcal{P}^*} \min_{\beta \in \mathcal{P}} l_2(\alpha, \beta), \quad (19)$$

Table: Comparison of transfer performance in same design space

Methodologies	A, B \rightarrow C		A, C \rightarrow B		B, C \rightarrow A	
	ADRS	HV	ADRS	HV	ADRS	HV
Ground Truth	0.0	0.0984	0.0	0.0684	0.0	0.0809
ANN-TL ⁵	0.069	0.0891	0.045	0.0643	0.031	0.0749
Deep-Ens ⁶	0.072	0.0840	0.066	0.0599	0.055	0.0727
ERM ⁷	0.080	0.0877	0.064	0.0629	0.043	0.0742
IRMv1 ⁸	0.025	0.0938	0.023	0.0667	0.027	0.0766
Ours	0.021	0.0944	0.017	0.0679	0.020	0.0796

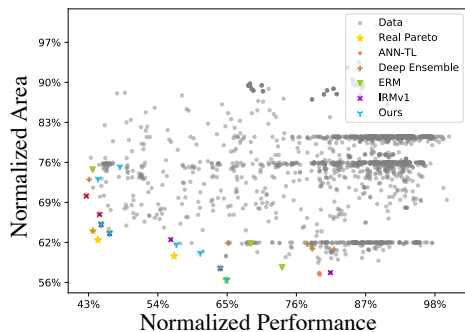
⁵Jianwang Zhai, Yici Cai, and Bei Yu (2023). “Microarchitecture Power Modeling via Artificial Neural Network and Transfer Learning”. In: *2023 28th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 1–6.

⁶Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell (2017). “Simple and scalable predictive uncertainty estimation using deep ensembles”. In: *Advances in neural information processing systems* 30.

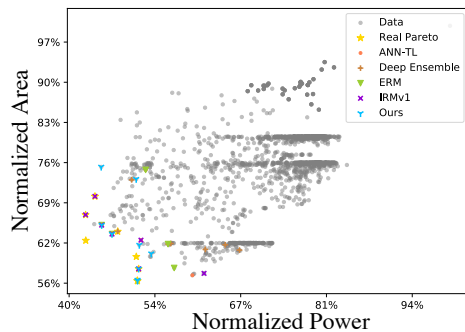
⁷VN Vapnik (1992). *Principles of risk minimization for learning theory*, *Advances in Neural Information Processing NIPS 4* (pp. 831±838).

⁸Martin Arjovsky et al. (2019). “Invariant risk minimization”. In: *arXiv preprint arXiv:1907.02893*. [27/30](#)

Pareto-optimal Sets Visulation



(a)



(b)

Pareto-optimal sets using source task data from same design space. Left: performance versus area; right: power versus area.

To the best of our best knowledge, all previous methods cannot spread knowledge across different design spaces.
This serves as an ablation study.

Table: Comparison of transfer performance in different design spaces

Methodologies	A, D \rightarrow C		A, C \rightarrow D	
	ADRS	HV	ADRS	HV
Ground Truth	0.0	0.0984	0.0	0.7792
w/o. Pre-train	0.0533	0.0857	0.0853	0.7465
w/o. Ensemble	0.0275	0.0892	0.0722	0.7531
w/o. IRM	0.0293	0.0910	0.0701	0.7569
w/. IRMv1	0.0232	0.0914	0.0687	0.7602
Ours	0.0217	0.0924	0.0641	0.7624

THANK YOU!