

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

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

¹The Chinese University of Hong Kong

²Huawei Noah's Ark Lab

³HiSilicon



Background and Motivation

- Microarchitecture design space exploration (DSE) determines the detailed structures of a microprocessor to obtain desired performance power and area (PPA) values.
- Complexity of design spaces and large evaluation cost are two major challenges.
- Transferring from historical explored tasks to new task reduces the exploration cost.
- Extract the relationship that is approximately invariant among different historical source design tasks, even equipped with different design spaces.

Typical Microprocessor Components

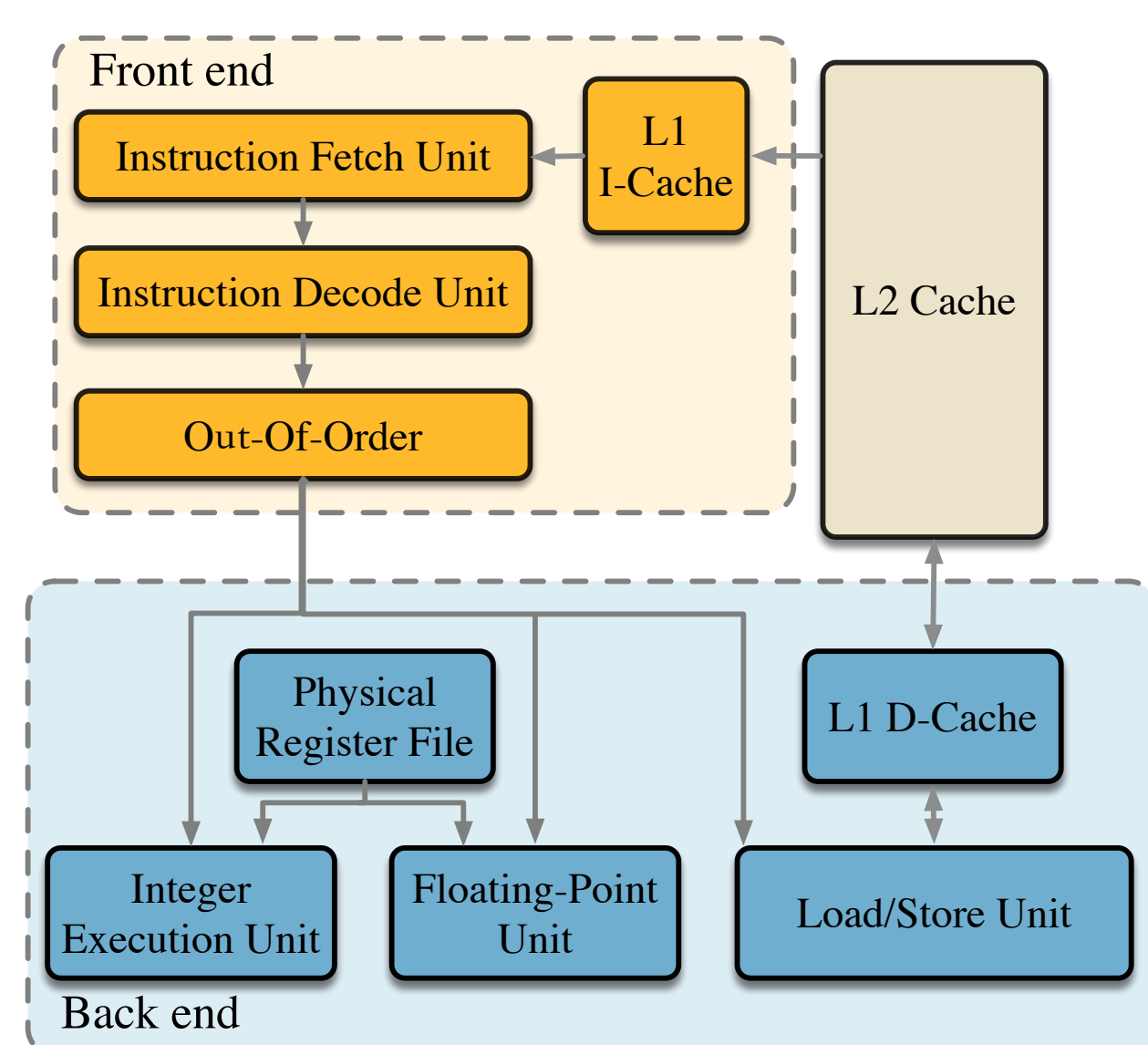


Figure 1. Visualization of typical microarchitecture pipeline.

Invariant Risk Minimization (IRM)

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:

$$\min_{\mathbf{u}, \mathbf{w}} \sum_{s=1}^S \mathcal{R}^s(\mathbf{u}, \mathbf{w}), \quad (1)$$

$$\text{s.t. } \mathbf{w} \in \operatorname{argmin}_{\mathbf{w}^s} \mathcal{R}^s(\mathbf{u}, \mathbf{w}^s),$$

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

$\mathcal{L}(\cdot)$: the loss function.

Simplified IRMv1:

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

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

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

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

Design Tasks

- Source task: dataset \mathcal{D}^s containing n_s parameter vectors $X^s = [\mathbf{x}_1^s, \dots, \mathbf{x}_{n_s}^s]^T \in \mathbb{R}^{n_s \times d_s}$ 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: contains only the set of legal parameter configuration vectors $X^t = [\mathbf{x}_1^t, \dots, \mathbf{x}_{n_t}^t]^T \in \mathbb{R}^{n_t \times d_t}$.

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 $X = \{\mathbf{x} | \mathbf{x} \succeq \mathbf{x}', \forall \mathbf{x}' \in X^t, Y = \{f(\mathbf{x}) | X \in X^t\}$.

Overall flow

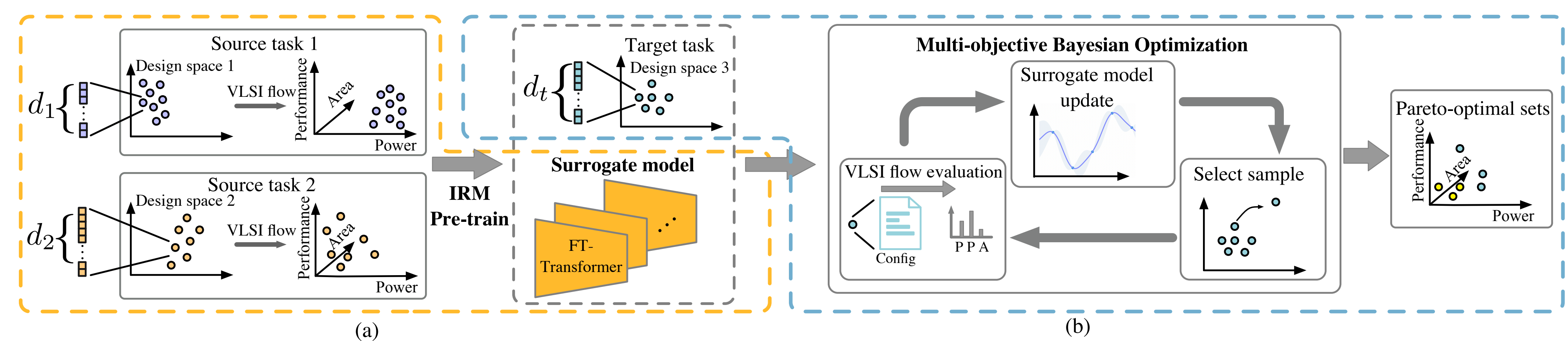


Figure 2. Workflow of proposed IT-DSE.

- 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

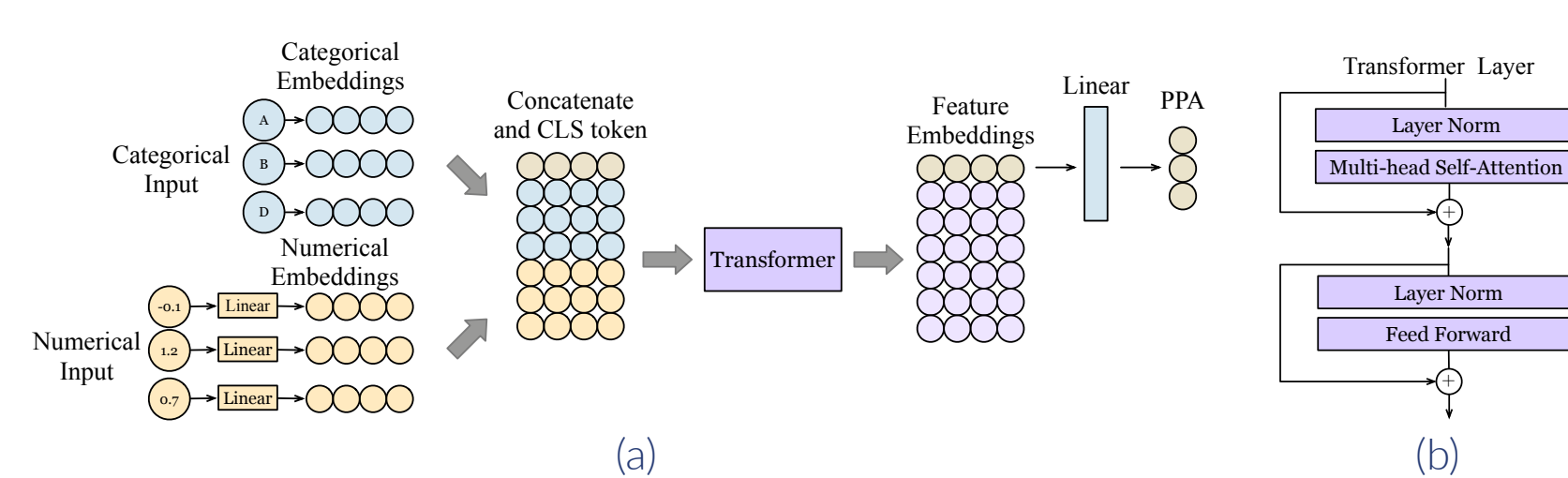


Figure 3. (a) FT-Transformer architecture. Firstly, Feature Tokenizer transforms input parameters into embeddings. The embedding is then processed by the Transformer module. (b) One Transformer layer.

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

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

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

Pre-training on Multiple Source Tasks:

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

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

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

Transformed embeddings:

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

Bayesian Invariant Risk Minimized Feature Extraction

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

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

Objective for Bayesian IRM:

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

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

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 $\phi_{\mathbf{u}}(\cdot)$.

Multi-objective Bayesian Optimization

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

$$p(\mathbf{y}_*^t | \mathbf{x}_*^t) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{y}_*^t | f'(X^t)}, \boldsymbol{\sigma}_{\mathbf{y}_*^t | f'(X^t)}^2), \quad (8)$$

- FT-Transformer Ensemble as Surrogate Model.
- Pareto EHVI Acquisition Function: $EHVI(\hat{\mathbf{x}}) = \mathbb{E}_{P(f'(\hat{\mathbf{x}}) | \mathcal{D})} (HVI(f'(\hat{\mathbf{x}}) | \mathcal{P}(\mathbf{y}), \mathbf{y}_{ref}))$

Transfer Design Space Exploration Experiment Setup

Table 1. 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.

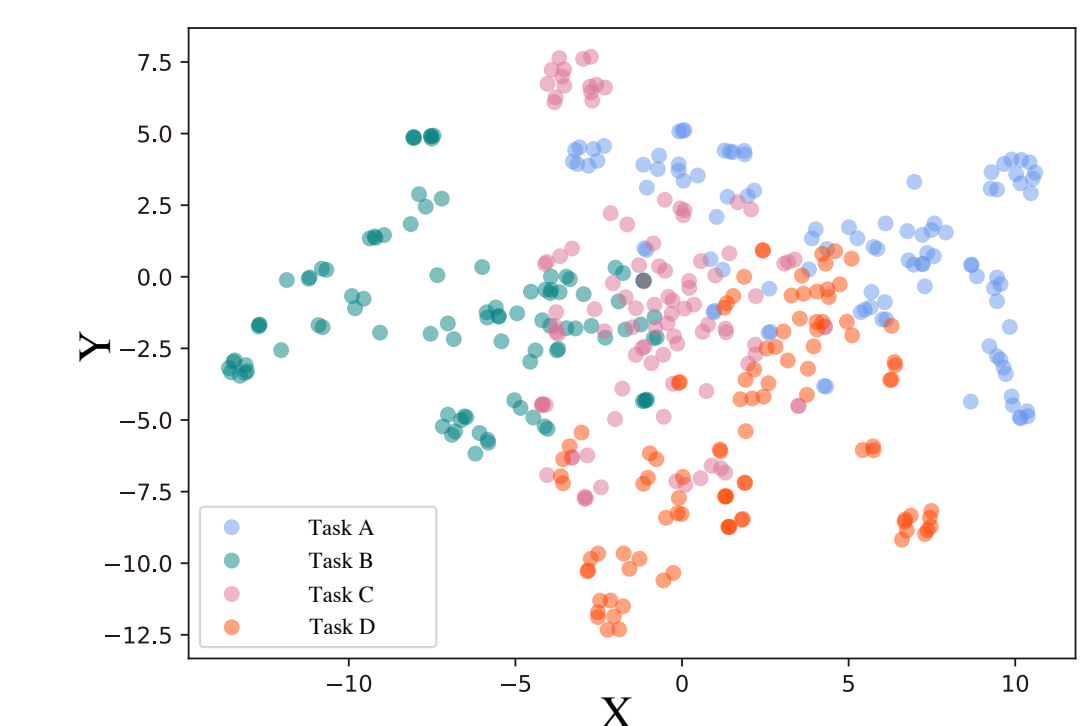


Figure 4. Data distribution for 4 tasks.

Evaluation Metrics:

Hypervolume:

$$HV_{\mathbf{y}_{ref}}(\mathcal{P}(\mathbf{y})) = \lambda_M(\cup_{\mathbf{y} \in \mathcal{P}(\mathbf{y})} [\mathbf{y}, \mathbf{y}_{ref}]), \quad (9)$$

Average distance to reference set (ADRS):

$$ADRS(\mathcal{P}^*, \mathcal{P}) = \frac{1}{|\mathcal{P}^*|} \sum_{\alpha \in \mathcal{P}^*} \min_{\beta \in \mathcal{P}} l_2(\alpha, \beta), \quad (10)$$

Transferring Performance

Table 2. Comparison of transfer performance in same design space

Methodologies	A, B → C		A, C → B		B, C → A	
	ADRS	HV	ADRS	HV	ADRS	HV
Ground Truth	0.0	0.0984	0.0	0.0684	0.0	0.0809
ANN-TL	0.069	0.0891	0.045	0.0643	0.031	0.0749
Deep-Ens	0.072	0.0840	0.066	0.0599	0.055	0.0727
ERM	0.060	0.0877	0.064	0.0629	0.043	0.0742
IRMv1	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

Table 3. Comparison of transfer performance in different design spaces

Methodologies	A, D → C		A, C → 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

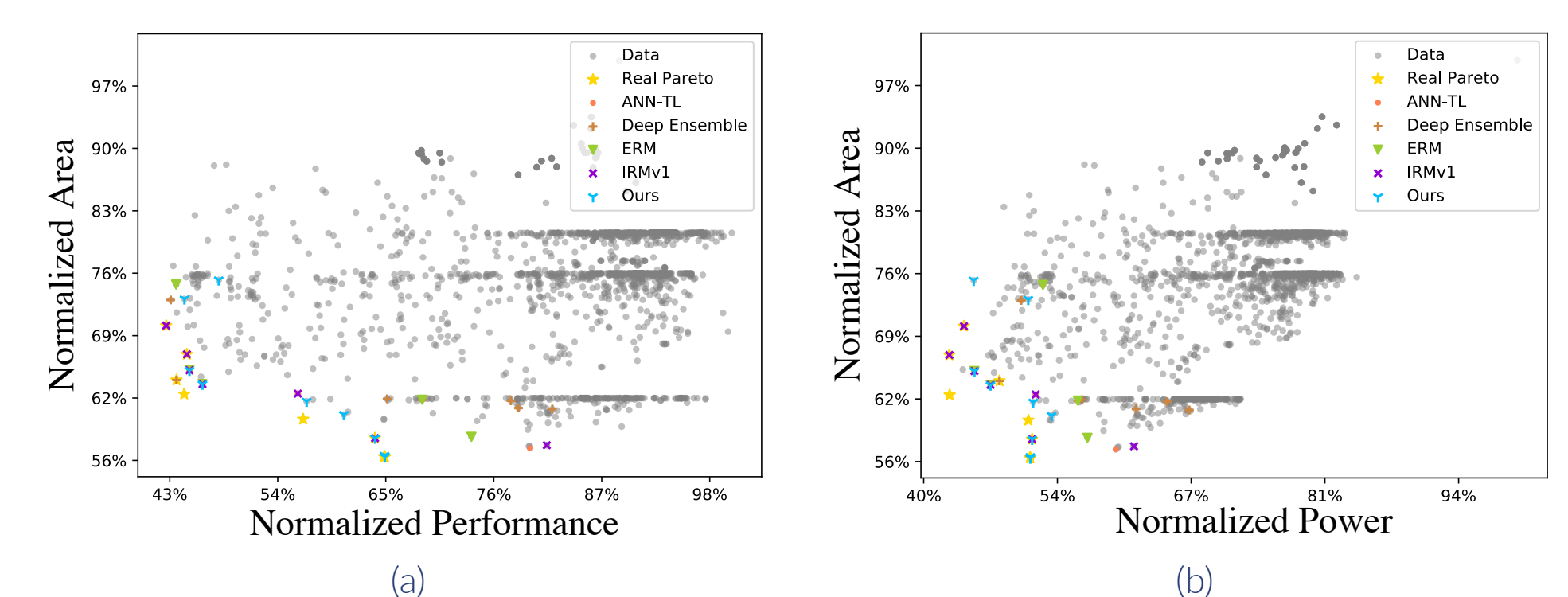


Figure 5. Pareto-optimal sets using source task data from same design space. Left: performance versus area; right: power versus area.