# [AI System Design – Term Project]
# Optimizations of
# AI Accelerator Networks

Chester Sungchung Park (박성정)

SoC Design Lab, Konkuk University

Webpage: http://soclab.konkuk.ac.kr

# Teaching Assistants

❑ Shinyoung Kim (shinyoungkim@konkuk.ac.kr)
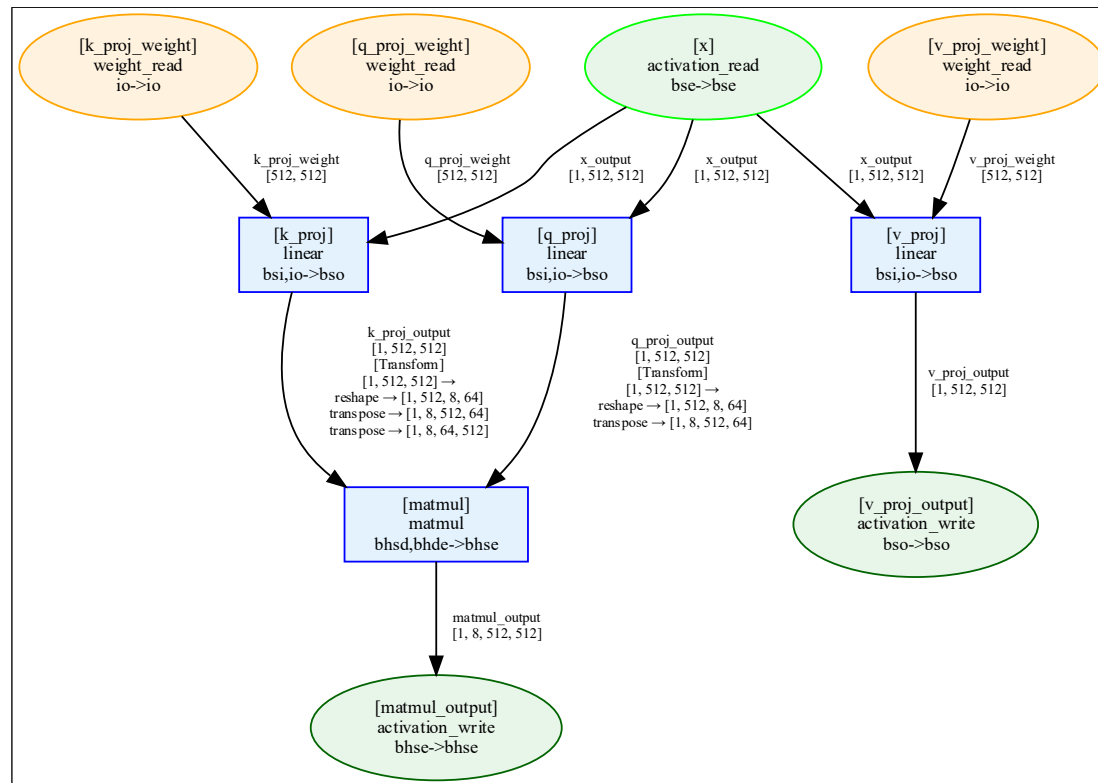
❑ Jaesuk Lee (jaesuki98@konkuk.ac.kr)

# Outline

❑ Introduction
- Target workload
- NetTLMSim

❑ Problem to solve

❑ Reference codes

❑ Evaluation

❑ Submission

# Target Workload

❑ BERT-medium *Layer group A*

- Layer group 0 + matmul (QK$^T$)
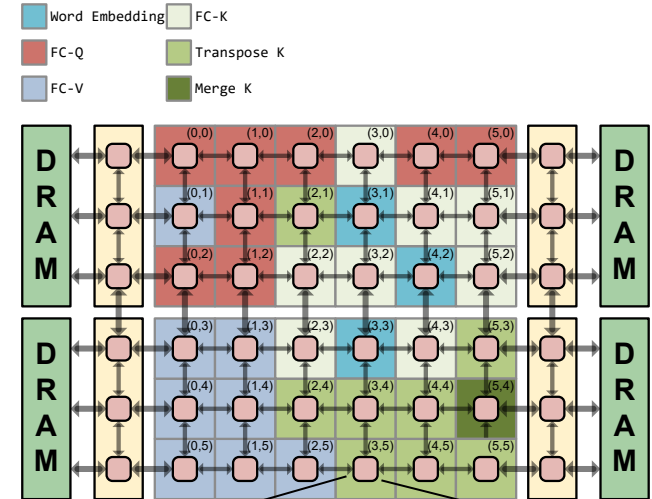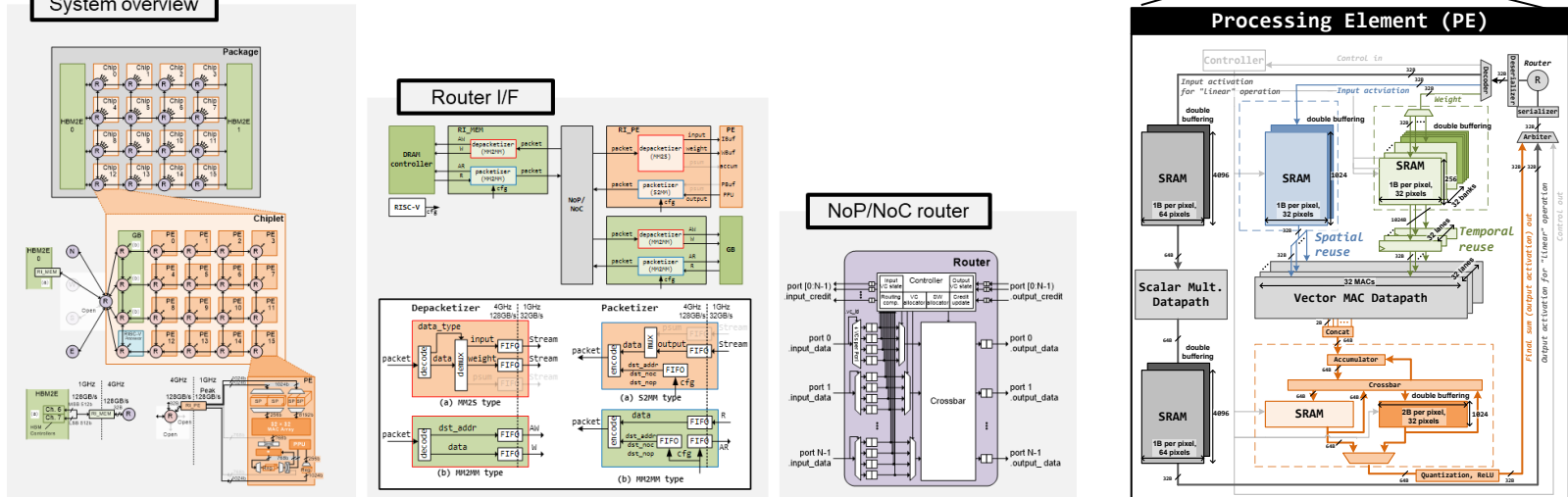
# NetTLMSim

## ☐ Virtual prototype

**Assumed design parameters**



| | | |
|---|---|---|
| **Space-Time Mapping** | Target Hardware | 16 chiplets per package [1], 16 PEs per chiplet [1, 2], 1024 MACs per PE [1, 3] Total 512 TOPS at 1 GHz freq. |
| | PE Style | NVDLA-like vector MAC array |
| | Loop Order | OS-LWS [4] |
| **Network Architecture** | Topology | Mesh [2] |
| | Routing | DOR YX [2] |
| | Flow Ctrl | Cut-through [2] |
| | Packet Len | 1 head flit, 16 body flits [2] |

| | | |
|---|---|---|
| **Memory Architecture** | DRAM Type | HBM2E, 2 devices Total 1024 GB/s |
| | SRAM Type | Dual port [3] |
| | Global Buf. Size, Bandwidth | 640 KB activation storage [4], 128 GB/s per NoC router |
| | Local Buf. Size, Bandwidth | 36 KB x2 (double buffered) [4], 128 GB/s at 1 GHz freq. |
| | Precision | 8 bits (24 bits for partial sums) [4] |

[1] Cai at el. 2024  [3] Keller at el. 2023
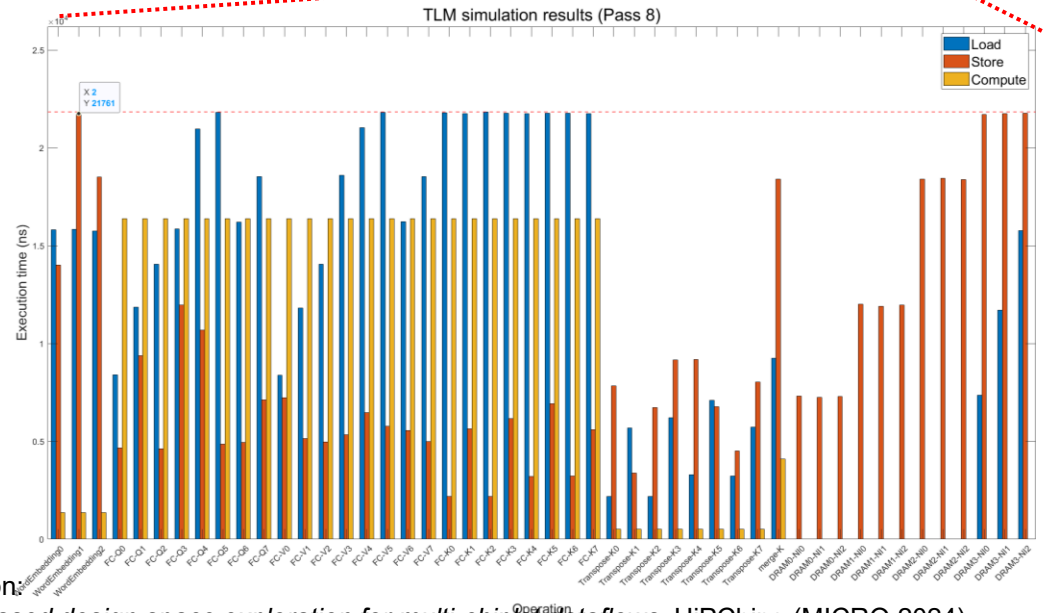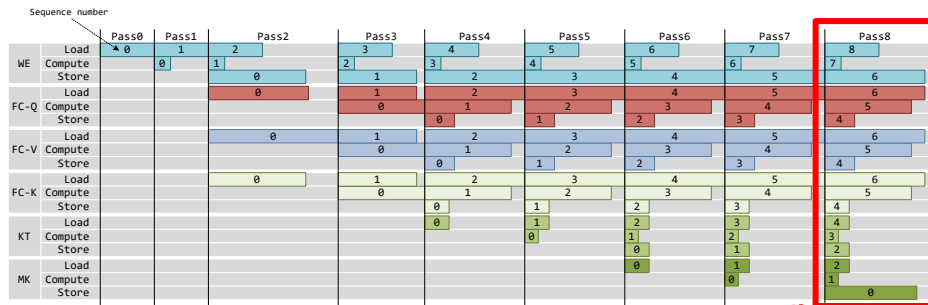[2] Zimmer at el. 2020  [4] Venkatesan at el. 2019

For more details, check the following presentation:
F. S. Park and C. S. Park, *Pre-RTL simulation based design space exploration for multi-chiplet dataflows*, HiPChips (MICRO 2024).
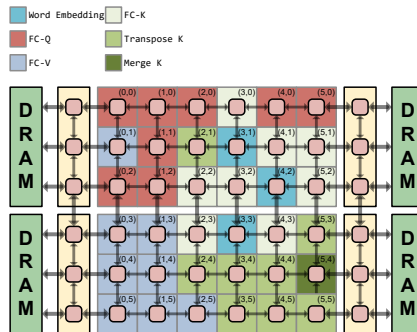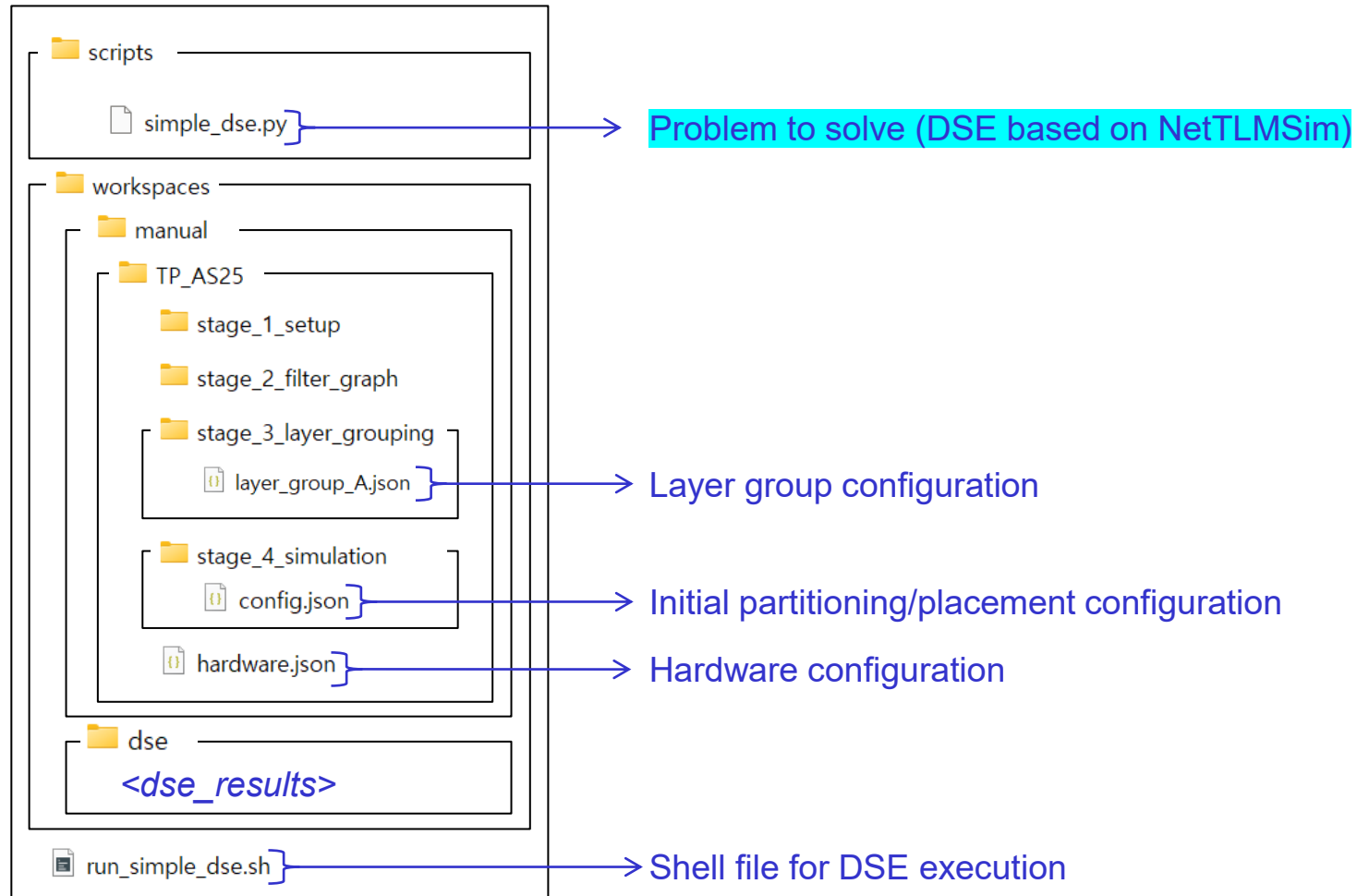
# NetTLMSim

## ❑ Virtual prototype (cont'd)



For more details, check the following presentation:
F. S. Park and C. S. Park, *Pre-RTL simulation based design space exploration for multi-chiplet dataflows*, HiPChips (MICRO 2024).
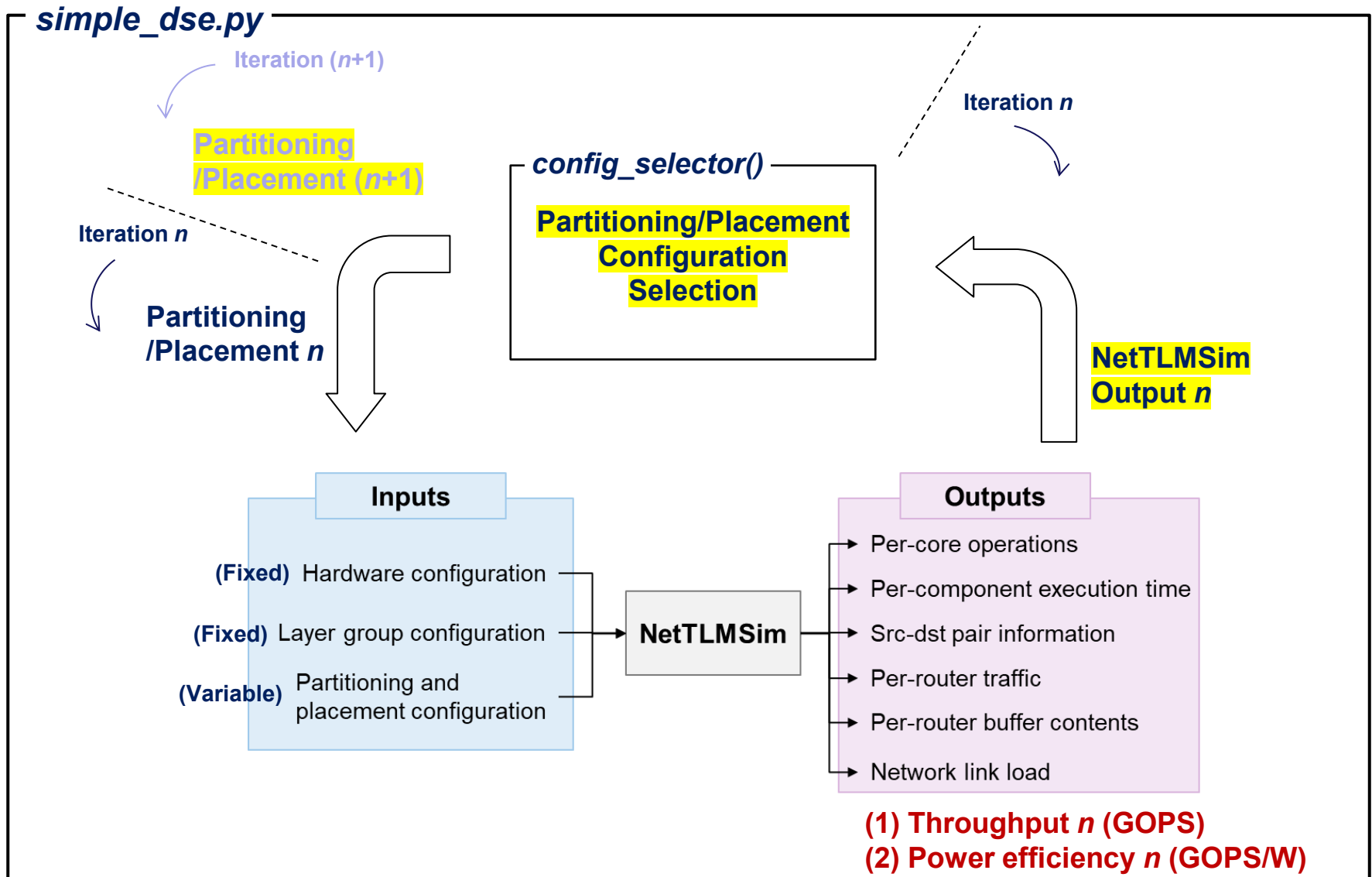
# Reference Project

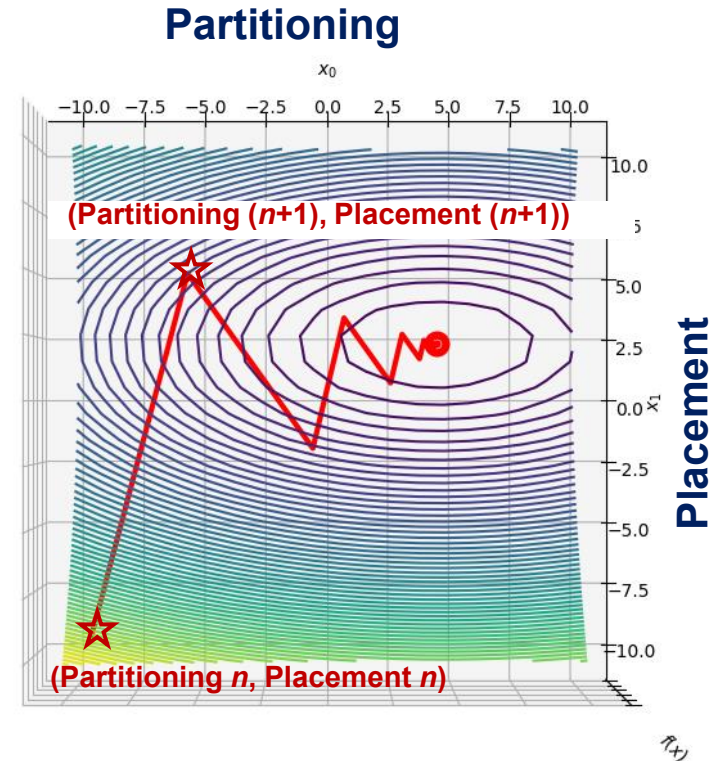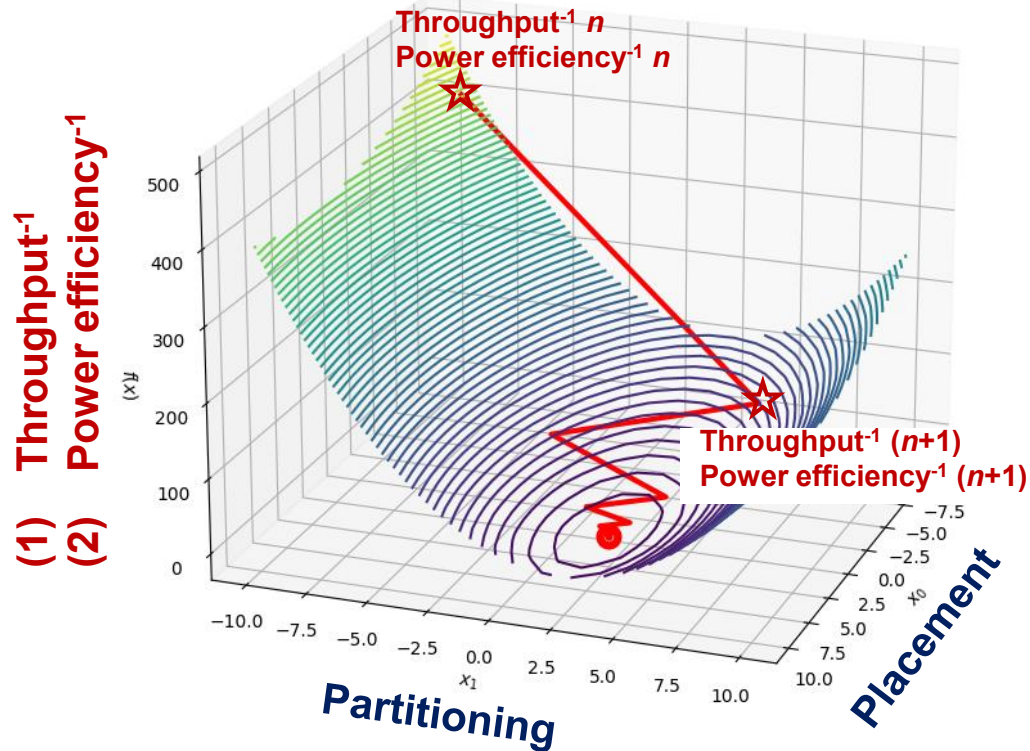❑ Build project and run simulation (See *Appendix A*)

```
📁 scripts
    📄 simple_dse.py  ─────────────→  Problem to solve (DSE based on NetTLMSim)
📁 workspaces
    📁 manual
        📁 TP_AS25
            📁 stage_1_setup
            📁 stage_2_filter_graph
            📁 stage_3_layer_grouping
                📄 layer_group_A.json  ───→  Layer group configuration
            📁 stage_4_simulation
                📄 config.json  ───→  Initial partitioning/placement configuration
            📄 hardware.json  ───→  Hardware configuration
    📁 dse
        <dse_results>
📄 run_simple_dse.sh  ───────────→  Shell file for DSE execution
```

# Problem to Solve

❑ Partitioning/placement configuration selection

- Input: **NetTLMSim output** for the *n*-th iteration
- Output: **NetTLMSim input** for the (*n*+1)-th iteration
  - ✓ Only partitioning/placement configuration *varied*
  - ✓ Hardware/layer group configuration *kept fixed*

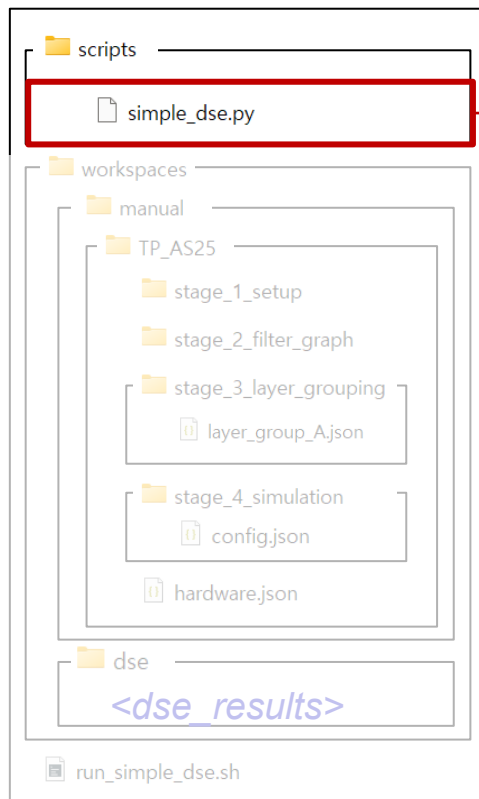**System-on-a-Chip Design LAB**

# Problem to Solve

# Problem to Solve

# Design Constraints

❑ Optimize with respect to partition and placement only
  - You are **not** allowed to modify any other settings, e.g., hardware, layer group configurations etc.
    - ✓ Otherwise, it won't be considered for evaluation

❑ Run *no* additional simulations (e.g., those using NetTLMSim) *inside* the partitioning/placement configuration selection
  - Only a single simulation run is allowed per iteration.

# Design Constraints

❑ Modify **only** (a part of) the body of *simple_dse.py*

- You are allowed to change **only** the paragraphs between *"Edit code below"* and *"Edit code above"*



Using these variables to utilize the results from the previous iteration

# Evaluation

❑ Submission completeness (10pt)

- Reproducibility

❑ Optimality (40pt)

- Final throughput and power efficiency

❑ DSE efficacy (30pt)

- Sampling efficiency
- General applicability

❑ Documentation (20pt)

- Covering all those mentioned in the above

# Submission Completeness

❑ Make sure that your submission is **complete**

- In other words, it should be possible to <span style="color:red">reproduce</span> your design together with the claimed throughput and power efficiency using **only** <u>the files that you submitted by the submission deadline</u>

# Optimality

❑ Evaluated both <span style="color:red">absolutely</span> and <span style="color:red">relatively</span>

- In other word, the quantity as well as the ranking matters

❑ Evaluate <span style="color:red">separately</span> for throughput and power efficiency

# DSE Efficacy

❑ Sampling efficiency

- Defined as the final quantity <span style="color:red">divided by</span> the total number of iterations
  - ✓ Evaluates *how many simulation runs are required to reach the final throughput and power efficiency*
- Considered for **both** *training* (fine-turning) and *inference* in the case of AI/ML (and separately, if necessary)
- The fewer iterations, the better sampling efficiency

❑ General applicability

- Evaluates whether the proposed idea is generally applicable
  - ✓ For example, different <span style="color:red">layer group configurations</span> (including layer groups 0~2)
- May be further verified by extra hardware/layer group configurations

# Submission

❑ Deadline
- **Dec. 12 (Fri)**, **10:00** GMT+9

❑ **Only one zip file** submission **per team** including the following files:
- Source code (all those needed to reproduce your results)
  - ✓ Except all those provided by TA (e.g., reference project)
- Documentation (PPT)
  - ✓ Including the explanation of the above source code

❑ Upload the zip file to the Ecampus
- Plus, send to chesterku2013@gmail.com as a *backup*

❑ You can post questions in the Ecampus (Q&A)

❑ Delayed submission will result in penalty!

# Appendix A:
# Running a DSE

# Runing a DSE

❑ Start the DSE

- Open the terminal with **'Ctrl + `'**

- Type the following command and press **'Enter'**
    - ✓ `./run_simple_dse.sh TP_AS25 --layer-groups A --iterations` **N**
    - ✓ Set iteration count to **N**

```
∨ TP_AS25_20251202_134201
  ∨ results
    ∨ iter_001
      > logs                        ⟶  Simulation outputs (Logs)
      {} config_changes.json         ⟶  Partition/placement configuration changes
      {} config.json                 ⟶  Partition/placement configuration at n-th iteration
      {} simulation_result.json      ⟶  Simulation outputs (Throughput, power efficiency, …)
      {} vis_report_config.json      ⟶  Visualization scripts
    > iter_002
```

```
  > iter_009
  > iter_010
  > stage_3_layer_grouping
  > stage_4_simulation            ⟶  Simulation inputs
  {} hardware.json
  {} run_log.json                 ⟶  Iteration history
  {} run_meta.json
  {} summary.json                 ⟶  DSE summary
```

System-on-a-Chip Design LAB

KU KONKUK UNIVERSITY