

Beyond Status-Quo NPUs: Hardware-Native Multi-Tenancy & Adaptive Inference

Stylianos I. Venieris

Sept 10th 2024
 Samsung Al Center, Cambridge, UK

Distributed Al Group @ Samsung Al



Stylianos (Stelios) Venieris



Rui Li







Alexandros Kouris



Hongxiang Fan

Konstantin Mishchenko

Hao (Mark) Chen

- Becoming ubiquitous across mobile devices, robots and cloud servers
- Backbone of deployed AI, enabling everything from facial recognition to advanced chatbots

٠



SAMSUNG Al Center

- Becoming ubiquitous across mobile devices, robots and cloud servers
- Backbone of deployed AI, enabling everything from facial recognition to advanced chatbots



SAMSUNG AI Center - Cambridge

- Becoming ubiquitous across mobile devices, robots and cloud servers
- Backbone of deployed AI, enabling everything from facial recognition to advanced chatbots

Existing NPUs are optimised for:

• CNNs | Single-DNN execution

New Challenges





SAMSUNG AI Center - Cambridge

Existing NPUs are optimised for:

• CNNs | Single-DNN execution

New Challenges







A New Class of NPUs: Multi-DNN Accelerators

- NPUs with native support for multi-model execution
- Key design decisions
 - Inter-DNN parallelisation strategy
 - DNN Engine Architecture
 - − Scheduling _____

Approach 1	Approach 2
Time-mux	Spatial Co-location
Enhancing existing architectures	Custom Designs
Static	Dynamic

SAMSUNG AI Center - Cambridge



S.I. Venieris, C.S. Bouganis, and N. D. Lane.

Multi-DNN Accelerators for Next-Generation AI Systems. IEEE Computer, 2023.

SAMSUNG AI Center - Cambridge

(a) Single-DNN Engine time-mux

Programmability

(b) Composable Engines time-mux & spatial co-location (c) Heterogeneous Engines(d) Dedicated Custom Enginestime-mux &time-mux &spatial co-locationspatial co-location



Customisation

S.I. Venieris, C.S. Bouganis, and N. D. Lane. Multi-DNN Accelerators for Next-Generation AI Systems. *IEEE Computer*, 2023.

SAMSUNG AI Center - Cambridge

(a) Single-DNN Engine time-mux

Programmability

(b) Composable Engines time-mux & spatial co-location (c) Heterogeneous Engines(d)time-mux &timspatial co-locationspate

S (d) Dedicated Custom Engines time-mux & spatial co-location



Customisation

S.I. Venieris, C.S. Bouganis, and N. D. Lane.

Multi-DNN Accelerators for Next-Generation AI Systems. IEEE Computer, 2023.

SAMSUNG Al Center - Cambridge



Customisation

S.I. Venieris, C.S. Bouganis, and N. D. Lane. Multi-DNN Accelerators for Next-Generation AI Systems. IEEE Computer, 2023.

Case Study: Sparse Multi-DNN Workloads

SAMSUNG Al Center - Cambridge

- Sparse multi-DNN workloads
 - Different sparsification methods, e.g. channel-wise, N:M block-wise, random
 - Diverse inference requests: <input, DNN, sparsification method, KPI>
- Key component: Scheduling
 - Determine the execution order of multi-DNNs
 - Focus on layer-wise execution on the NPU



Input user requests

Case Study: Sparse Multi-DNN Workloads

- Current approaches
 - Lack of adequate sparsity support
 - Only consider average sparsity rate
- Untapped opportunity: Fine-grained sparsity information
 - Sparsity pattern
 - Sparsity dynamicity
- Sparsity pattern \rightarrow hardware efficiency variability
 - NVIDIA Tensor Core: N:M block-wise sparsity
 - Samsung Exynos NPU: Random sparsity
- Sparsity dynamicity \rightarrow Runtime variance
 - Varied latency performance can affect the best schedule



Case Study: Dysta

- **Problem:** Schedule variably sparse DNNs on zero-skipping NPUs, when each request has different <input, DNN, sparsity pattern, runtime sparsity level, KPI>
- Approach: Sparsity-aware bi-level scheduling
 - Capture both sparsity pattern (static) and dynamic sparsity (run-time) information
 - SW-based static scheduler: Initial schedule
 - HW-based dynamic scheduler: Run-time schedule refinement



Case Study: Dysta

- **Problem:** Schedule variably sparse DNNs on zero-skipping NPUs, when each request has different <input, DNN, sparsity pattern, runtime sparsity level, KPI>
- Approach: Sparsity-aware bi-level scheduling
 - Capture both sparsity pattern (static) and dynamic sparsity level (run-time) information
 - SW-based static scheduler: Initial schedule-
 - HW-based dynamic scheduler: Run-time schedule refinement



Case Study: Dysta

- SAMSUNG AI Center - Cambridge
- **Problem:** Schedule variably sparse DNNs on zero-skipping NPUs, when each request has different <input, DNN, sparsity pattern, runtime sparsity level, KPI>
- Approach: Sparsity-aware bi-level scheduling
 - Capture both sparsity pattern (static) and dynamic sparsity level (run-time) information



- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm



- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm



- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm



- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm



- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm





SAMSUNG Al Center

- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm





SAMSUNG AI Center

- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm



SAMSUNG Al Center

Compute Unit

Track runtime

sparsity

- Hardware scheduler
 - FIFOs: Track per-task information
 - Monitor: Track runtime sparsity
 - Compute unit: Estimate latency and update schedule
 - Reconfigurable design: Reuse resources for latency estimation and scheduling algorithm
- Up to 10% lower latency constraint violation rate and nearly 4x lower average normalised turnaround time across request traffic levels over state-of-the-art methods





SAMSUNG Al Center

- Foundation models
 - Powerful pretrained models
 - Parameter-efficient Fine-tuning (PEFT) to downstream tasks
 - Frozen shared part & task-specific adaptation part



SAMSUNG AI Center - Cambridge

- Foundation models
 - Powerful pretrained models
 - Parameter-efficient Fine-tuning (PEFT) to downstream tasks
 - Frozen shared part & task-specific adaptation part



[2] L. Chen, et al. Punica: Multi-Tenant LoRA Serving. *MLSys*, 2024.
[3] Sheng, Ying, et al.S-LoRA: Serving Thousands of Concurrent LoRA Adapters. *MLSys*, 2024.

- Foundation models
 - Powerful pretrained models
 - Parameter-efficient Fine-tuning (PEFT) to downstream tasks
 - Frozen shared part & task-specific adaptation part



[2] L. Chen, et al. Punica: Multi-Tenant LoRA Serving. *MLSys*, 2024.[3] Sheng, Ying, et al.S-LoRA: Serving Thousands of Concurrent LoRA Adapters. *MLSys*, 2024.

- Foundation models
 - Powerful pretrained models
 - Parameter-efficient Fine-tuning (PEFT) to downstream tasks
 - Frozen shared part & task-specific adaptation part



[2] L. Chen, et al. Punica: Multi-Tenant LoRA Serving. *MLSys*, 2024.[3] Sheng, Ying, et al.S-LoRA: Serving Thousands of Concurrent LoRA Adapters. *MLSys*, 2024.

SAMSUNG Al Center

ons

SAMSUNG Al Center

- Cambridge

- Foundation models
 - Powerful pretrained models
 - Parameter-efficient Fine-tuning (PEFT) to downstream tasks
 - Frozen shared part & task-specific adaptation part
- Research Questions
 - How to swap efficiently between Adaptation layers? Different memory management approach?
 - How to design the underlying hardware?
 Specialised units or unified for Shared and Adaptation parts?
 - Model-Hardware Co-Design for the next milestone



[2] L. Chen, et al. Punica: Multi-Tenant LoRA Serving. *MLSys*, 2024.
 [3] Sheng, Ying, et al.S-LoRA: Serving Thousands of Concurrent LoRA Adapters. *MLSys*, 2024.

Thank you

Multi-DNN Accelerators for Next-Generation AI Systems

Stylianos I. Venieris Samsung Al

Chistos-Savvas Bouganis Imperial College London

Nicholas D. Lane University of Cambridge Samsung Al



Sparse-DySta: Sparsity-Aware Dynamic and Static Scheduling for Sparse Multi-DNN Workloads

Hongxiang Fan hongxiangfan@ieee.org Samsung AI Center & University of Cambridge Cambridge, UK

n Stylianos I. Venieris* org s.venieris@samsung.com & Samsung AI Center idge Cambridge, UK Alexandros Kouris a.kouris@samsung.com Samsung AI Center Cambridge, UK Nicholas D. Lane ndl32@cam.ac.uk University of Cambridge & Flower Labs Cambridge, UK



Stylianos I. Venieris

Cambridge

Samsung Al