

打破交互墙：以DLPU为中心的深度学习计算系统 Breaking the Interaction Wall: A DLPU-centric Deep Learning Computing System

杜子东¹, 郭崎¹, 赵永威¹, 曾惜¹, 李玲², 程丽敏² 等

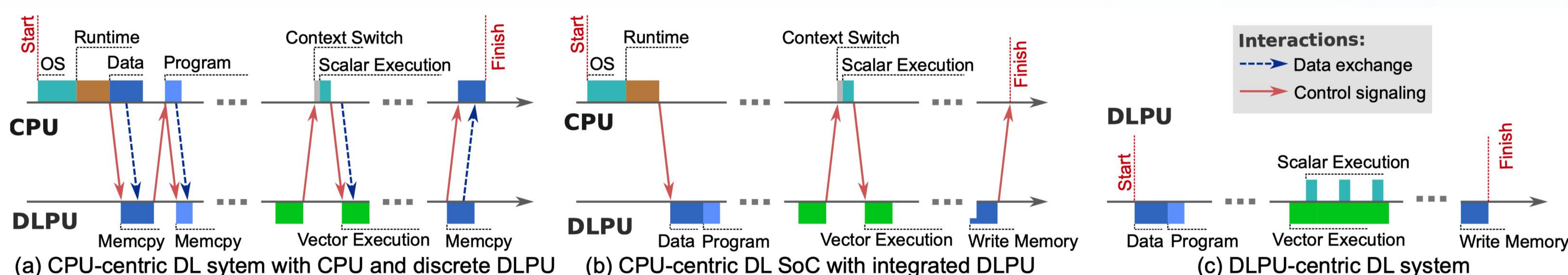
¹中国科学院计算技术研究所 ²智能软件研究中心, 中国科学院软件研究所

期刊: IEEE Transactions on Computers, 2020

联系人: 李玲 联系方式: 18500300237 liling@iscas.ac.cn

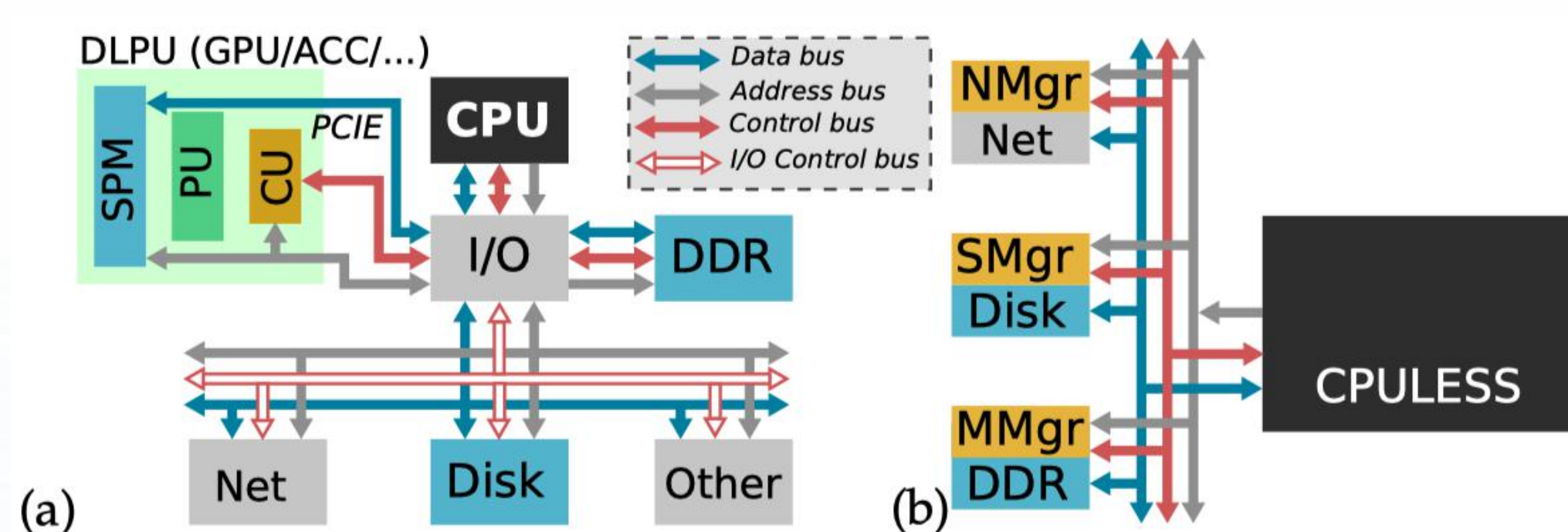
动机

- 以CPU为中心的深度学习计算系统存在交互墙问题, 导致系统效率低下
- 提出以DLPU (深度学习处理单元) 为中心的深度学习计算系统, 包括面向异常的编程模型 (EOP) 和CPULESS DLPU架构支持, 可以显著提升速度、降低能耗



DLPU为中心的深度学习计算系统

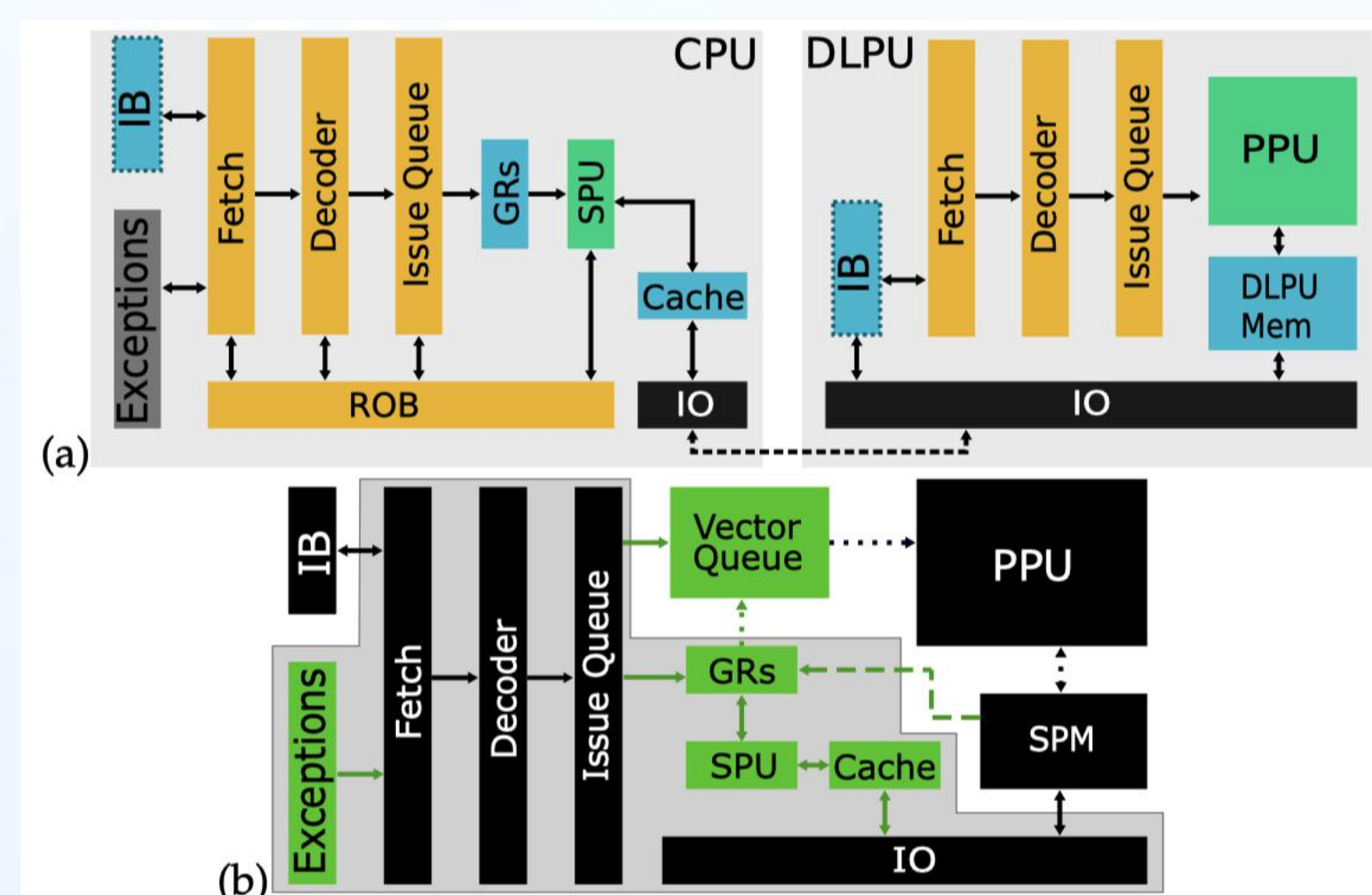
- 深度学习计算系统架构



(a) 主流的CPU为心的计算系统架构

(b) 以DLPU为心的计算系统架构

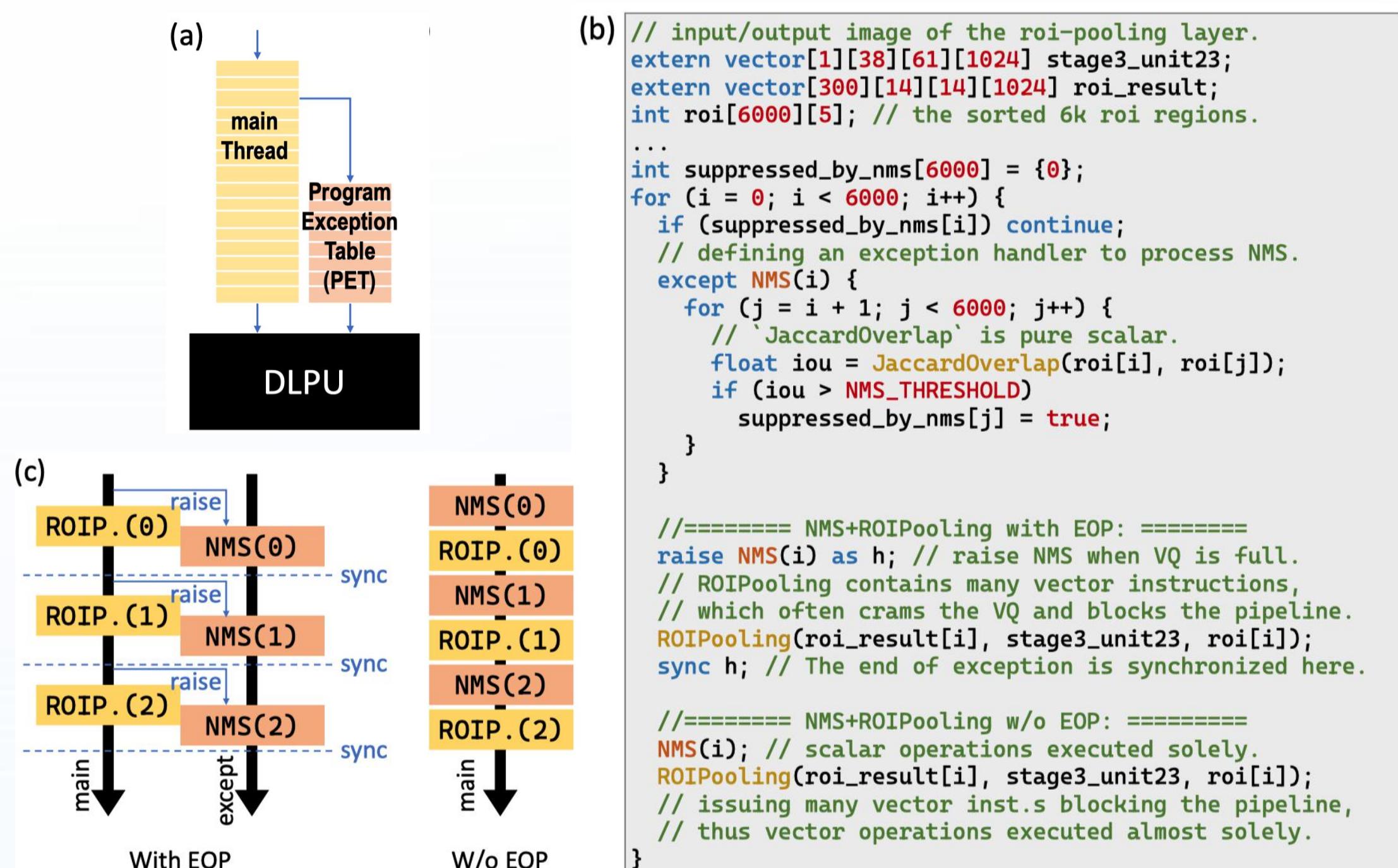
- CPULESS微架构



(a) CPU为心的微架构

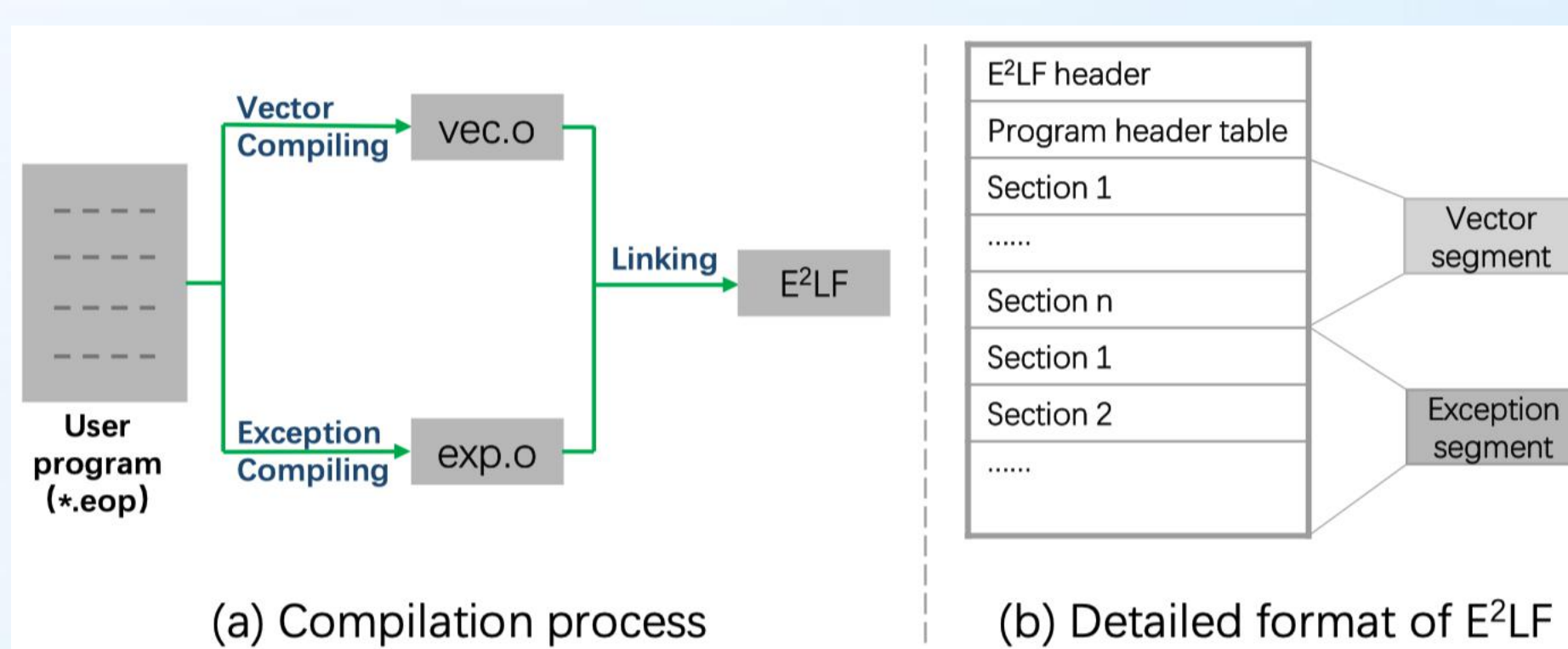
(b) CPULESS 微架构

- EOP编程模型



(a) 抽象机器模型 (b) Faster R-CNN算法的示例代码 (c) 程序员视角

- 编译链接

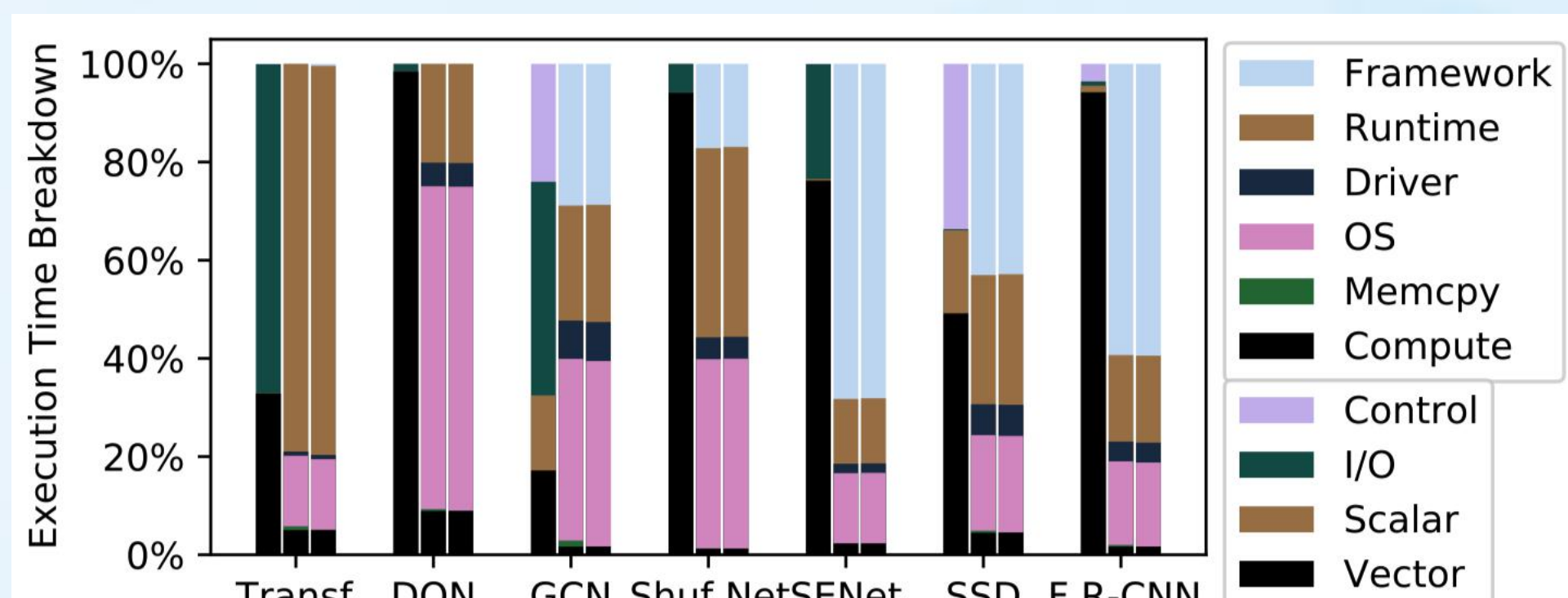
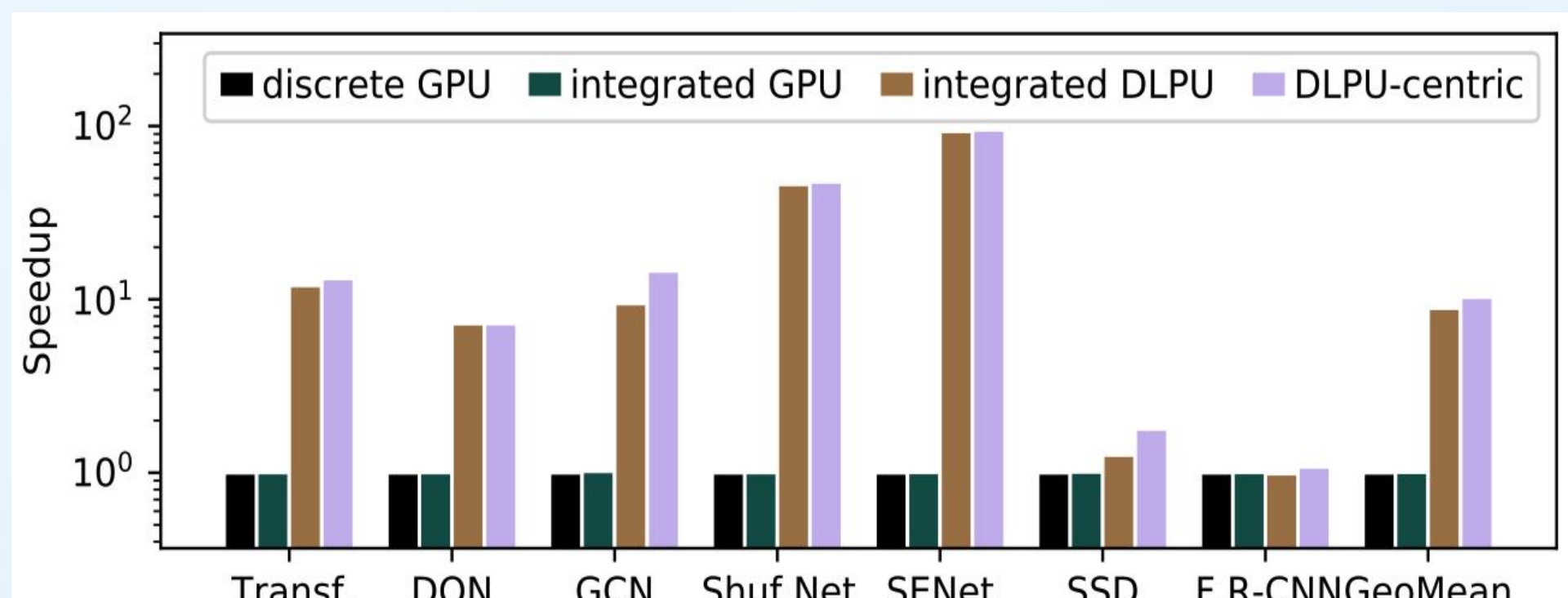


(a) Compilation process

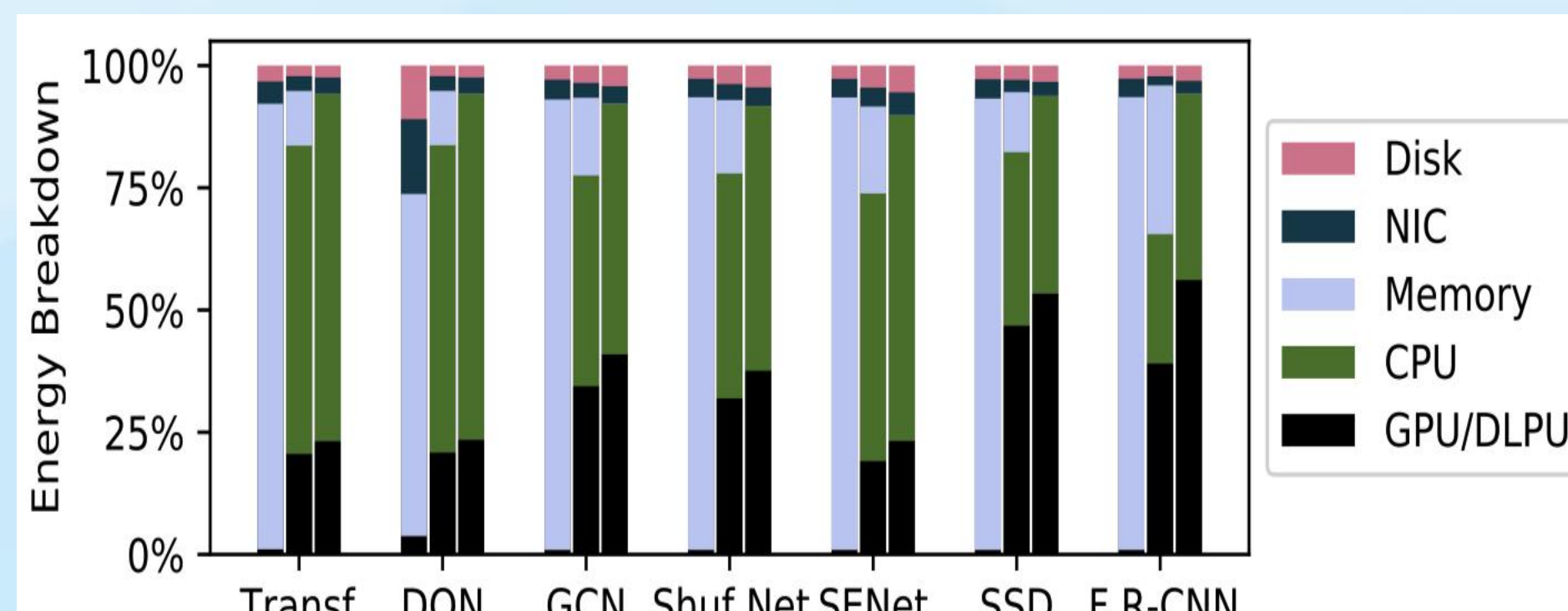
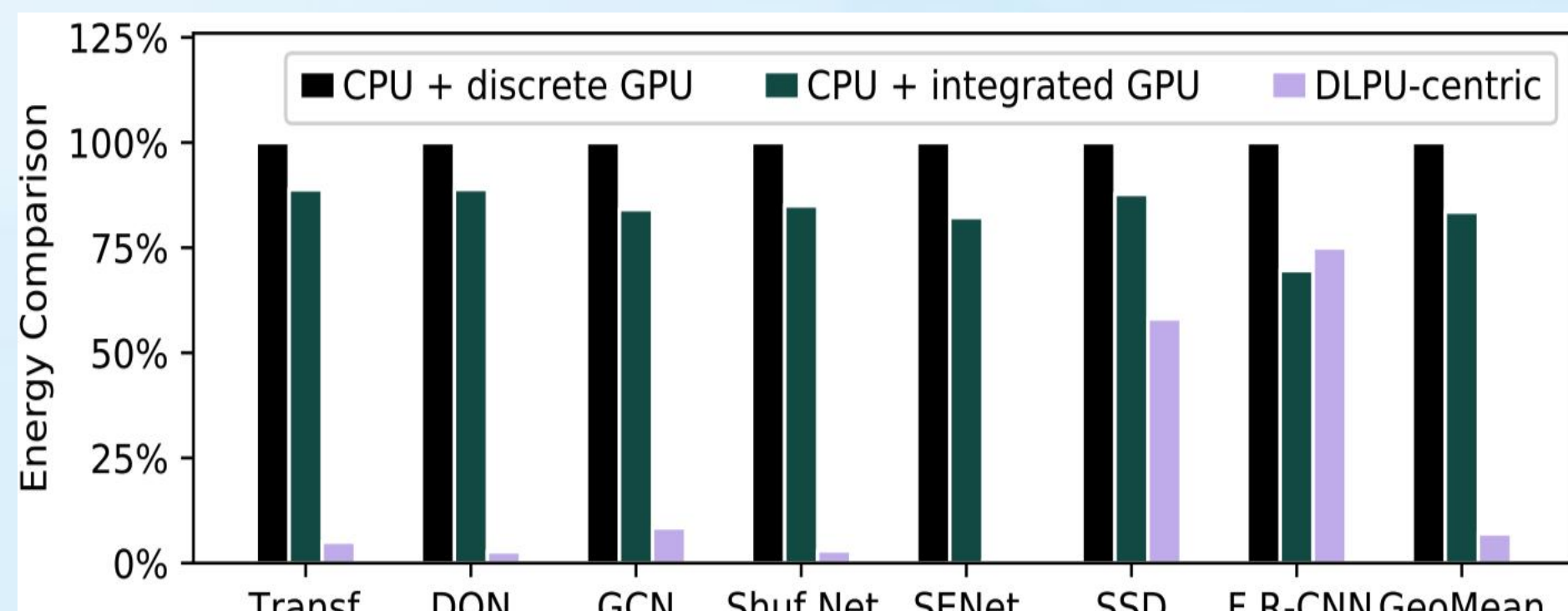
(b) Detailed format of E2LF

实验结果

- 对比硬件: Intel Xeon 6130 CPU, Nvidia Tesla V100 GPU
- 加速比



- 能耗



相对于主流的CPU+V100系统, 本文系统平均加速10.30倍, 节省92.99%的能耗