

From DSLs to Accelerator-Rich Platform Implementations: Addressing the Mapping Gap

Presenters: Steven Lyubomirsky*, Thierry Tambe*

Bo-Yuan Huang* Yi Li Mike He Gus Smith

Gu-Yeon Wei Aarti Gupta Sharad Malik

Zachary Tatlock

*Equal contribution



Harvard John A. Paulson School of Engineering and Applied Sciences PRINCETON

School of Engineering and Applied Science



The Mapping Gap





High-level application logic

Varied, low-level interfaces

FlexNLPU (Tambe *et al.*, 2020)



FlexNLPU (Tambe *et al.*, 2020)



FlexNLPU (Tambe *et al.*, 2020)



FlexNLP (Tambe *et al.*, 2020)



Compiler Support Challenges

- Lack of a unified hardware specification
 - No ISA!
 - Interfaces via MMIO commands
- Granularity mismatch
 - Source language is fine-grained
 - Accelerators operations are high-level
- Custom memory and numerical formats

Instruction-Level Abstraction (Huang *et al.*, 2019)

MMIO Access	Instruction		Αι	rch
WR 0x33000010	GB_CONTROL_START		(0	:fg
WR 0x33000050	GB_ATTENTION_START	ſ	-	h
WR 0x33700010	CONFIG_GB_CONTROL	Н		δ
WR 0x33700050	CONFIG_GB_ATTENTION			
•••	•••			

Architectural state variable:

(cfgMode, cfgIsRnn, cfgNumTimestep, valid, memL, ...)

Instruction CONFIG_GB_CONTROL:

$$\begin{split} &\delta_i \stackrel{\text{def}}{=} (inWr \land \neg inRd \land inAddr = 0x33700010) \\ &N_i[cfgMode] \stackrel{\text{def}}{=} ITE(valid, inData[1:3], cfgMode) \\ &N_i[cfgIsRnn] \stackrel{\text{def}}{=} ITE(valid, inData[4:4], cfgIsRnn) \end{split}$$

Instruction-Level Abstraction (Huang *et al.*, 2019)

Formally specified interface!				
Just like an ISA!) 1			

WR 0x33000050	GB_ATTENTION_START
WR 0x33700010	CONFIG_GB_CONTROL
WR 0x33700050	CONFIG_GB_ATTENTION

. . .

. . .

Architectural state variable:

(cfgMode, cfgIsRnn, cfgNumTimestep, valid, memL, ...)

```
Instruction CONFIG_GB_CONTROL:
```

```
\begin{split} &\delta_i \stackrel{\text{def}}{=} (inWr \land \neg inRd \land inAddr = 0x33700010) \\ &N_i[cfgMode] \stackrel{\text{def}}{=} ITE(valid, inData[1:3], cfgMode) \\ &N_i[cfgIsRnn] \stackrel{\text{def}}{=} ITE(valid, inData[4:4], cfgIsRnn) \end{split}
```

Instruction-Level Abstraction (Huang *et al.*, 2019)

MMIO Access	Instruction	to MIMIO commands!		
WR 0x33000010	GB_CONTROL_START	(сј дмоае, сј дізкпп, сј дмити imestep, valid, memL,)		
WR 0x33000050	GB_ATTENTION_START	Instruction CONFIG_GB_CONTROL:		
WR 0x33700010	CONFIG_GB_CONTROL	$\delta_i \stackrel{\text{def}}{=} (inWr \wedge \neg inRd \wedge inAddr = 0x33700010)$		
WR 0x33700050	CONFIG_GB_ATTENTION	$N_{i}[cfgMode] \stackrel{\text{def}}{=} ITE(valid, inData[1:3], cfgMode)$ $N_{i}[cfgIsBnn] \stackrel{\text{def}}{=} ITE(valid, inData[4:4], cfgIsBnn)$		
•••		$n_i[c] giskinj = n L (valia, inDala[4, 4], c) giskin)$		

Easy loworing



TVM Bring Your Own Codegen



Syntactic pattern matching is key!

3LA: Codegen via ILA

Import from frameworks to TVM Match Relay fragments



Map Relay fragments to ILA instructions

FlexNLP func A: instr. 1 (STR r2, 0xffff0000) instr. 2 (LDR r3, 0xffff0010)

FlexNLP func B: instr. 3 (STR r2, 0xffffaa00) instr. 4 (LDR r3, 0xffffaabb)

NVDLA func C: instr. 5 (STR r4, 0xffff0100) instr. 6 (LDR r5, 0xffff0110)



Lower to accelerator interface

; CPU instructions CMP r0, r1 SUBGT r0, r0, r1 BNE loop ; Access accel 1 (MMIO) r2, 0xffff0000 STR r3, 0xffff0010 LDR ; Access accel 2 (MMIO) r4, 0xffff0100 STR r5, 0xffff0110 LDR ; CPU instructions r3, r2 MOV SUBGT r0, r0, r1 lr В

Heterogeneous hardware backends



Further Possibilities

- End-to-end functional verification
- Flexible matching (semantic reasoning)
- Automating compiler stack generation

Thank You!



Bo-Yuan Huang



ang Thierry Tambe



Yi Li



Mike He



Gus Smith



Gu-Yeon Wei



Aarti Gupta



Sharad Malik



Zachary Tatlock

Thank You!





Joint University Microelectronics Program

www.src.org/program/jump



Semiconductor Research Corporation

@srcJUMP



TVM (Chen et al., 2018)

Itvm

Open source, automated, optimization framework for deep learning.

