Machine learning and Optimization-driven Compilers for Heterogeneous Architectures: MOCHA

Howard Shrobe DARPA Information Innovation Office (I2O)





Machine learning and Optimization-guided Compilers for Heterogeneous Architectures (MOCHA)

Compilers can be rearchitected to generate highly efficient code for heterogeneous hardware platforms by employing ML-generated components that can be incorporated with minimal human effort



Today's compiler architecture

MOCHA

Apply ML across the entire compiler pipeline

- Support multiple compute architectures simultaneously
- Enable holistic optimization across passes

IR: Intermediate Representation CPU: Central Processing Unit GPU: Graphical Processing Unit TPU: Tensor Processing Unit



The tale of the F-22 "Software Crisis"



"As of 2003, after over a decade of effort, the code was crashing, on average, every three hours. Some parts of the software system were failing every 90 minutes."

SLOC: Source Lines Of Code

https://freerepublic.com/focus/f-news/1110502/posts



The lesson of F-22 software

- DoD system design, development and deployments happen over very long timeframes
- Processors come and go much more quickly
 - Even new versions of existing components may have wildly different micro-architecture with different performance profiles
- Mission needs change much more quickly
- Writing software "close to the metal" to gain performance impedes agility and adaptability
- Mission software should be "hardware agnostic" and written at a high level of abstraction
- The **compiler should decide** how to utilize different processors and accelerators, enabling the software to be remapped to take advantage of new hardware
 - Increases agility
 - Requires much less (scarce) human resources
 - Increases software reliability and correctness



Today: State-of-the-art architectures include heterogeneity

End of Moore's Law has led to proliferation of processing units and created a need for specialized computer architectures and accelerators



Open Multimedia Applications Platform system on chip: TI OMAP4 SoC

- More computational power on accelerators than on CPU cores
- Specialization overcomes power limitations
- Programming heterogenous systems is difficult



2024 and beyond



With Compiler Support, Heterogeneous and Specialized SoC Outperforms CPU's and GPU's



Transmuter v.5 SoC (DARPA MTO SDH Program UMich Team) Energy Efficiency for Infrared Small Target Detection



Source: DARPA SDH program

Manual development of compiler tool chain required

SoC: System on Chip SDH: Software-Defined Hardware SWaP: Size, Weight, and Power

Pal, Subhankar, et al. "Transmuter: Bridging the efficiency gap using memory and dataflow reconfiguration." *Proceedings of the ACM International Conference on Parallel Architectures and Compilation Techniques*. 2020.



Today: Compiler development is dependent on manual effort



- Unresponsive to hardware innovations
- Abstraction of the standard CPU architecture
- Based on classic Intermediate Representations (IR), analyses, and transformations; ad hoc heuristics; and hand-tuned optimizations
- Manufactures spec-sheets are inaccurate due to microarchitectural interactions
- Don't match the capabilities of accelerators
- Optimization design is a "black art"



Source: Chris Cummins, Deep Learning for Compilers 2019

DISTRIBUTION STATEMENT A: Approved for public release: distribution is unlimited.

Compiler Performance





CPU: Central Processing Unit DSP: Digital Signal Processor GPU: Graphical Processing Unit TPU: Tensor Processing Unit



Tomorrow: Machine learning can break the manual bottlenecks

Technical hypothesis: Compilers can be rearchitected to generate highly efficient code for heterogeneous hardware platforms by using ML-generated components that can be incorporated with minimal human effort



Approach:

- Move to "Hardware Agnostic" programming
- Reduce human effort in compiler construction by using Machine Learning to generate compiler components
- Target a variety of accelerators, producing correct, highly optimized code for a heterogenous ensemble of components



- **Front-end** Front-end components need to identify regions of the code that can be targeted to accelerators
 - Candidate regions identified by examining the call graph, which is constructed in the compiler front end
 - Performance models identify which sub-graphs are candidates for acceleration



LUT: Logical Unit



Speedup of H.264 decoder by implementing identified hardware accelerators

Limitations: Performance models are ad hoc and not general

Zacharopoulos, Georgios, et al. "Compiler-assisted selection of hardware acceleration candidates from application source code." *2019 IEEE 37th International Conference on Computer Design (ICCD)*. IEEE, 2019.



Reason for confidence: VeGen: Learning what optimization to apply



Support for Programming Languages and Operating Systems. 2021.



A matrix-vector multiplication kernel from TVM

ICC: Intel C Compiler	
Clana: LLVM compiler front-end	

Clang: LLVM (compiler front-er
TVM: Tensor	Virtual Machine
VeGen: Vecto	orizer Generator

void		ICC (1x)	GCC (1.5x)	Clang/LLVM (2.2)
dot 16x1x16 uint8 int8 int32(<u>s</u> , (
uint8 t data[restrict 4].	_			
int the real [restrict 16][/]	movzx	rlld, [rdi]	vmovdqa xmm0, .LC0[rip]	movzx eax, [rdi + 3]
IIICO_U KEIHEI[FESUFICU 10][4],	movsx	eax, [rsi]	vmovdqu xmm1, [rsi]	vpbroadcastd zmm8, eax
<pre>int32_t output[restrict 16]) {</pre>	imul	rlld, eax	vmovdqu xmm3, [rsi+16]	movzx eax, [rdi + 2]
for (int i = 0; i < 16; i++)	movsx	ecx, [1+rs1]	vmovdqu xmm2, [rsi+32]	vpbroadcastd zmm9, eax
$\mathbf{f}_{\mathbf{r}} = (\mathbf{r}_{\mathbf{r}} + \mathbf{r}_{\mathbf{r}}) + (\mathbf{r}_{\mathbf{r}} + \mathbf{r}_$	movzx	r8d, [1+rd1]	vmovdqu xmm6, [rsi+48]	movzx eax, [rdi + 1]
IOF (INC $K = 0; K < 4; K^{++}$)	imul	rød, ecx	vpand xmm10, xmm0, xmm1	vpbroadcastd zmm10, eax
output[i] +=	movzx	rlud, [2+rdi]	vpsrlw xmml, xmml, 8	movzx eax, [rd1]
data[k] * kernel[i][k];	movsx	r9d, [2+rs1]	movzx r9d, [rd1]	vpbroadcastd zmmll, eax
1	imul	riud, ryd	vpand xmm4, xmm0, xmm3	vmovdqu xmmU, [rsi]
}	movsx	eax, [3+rs1]	vpsrlw xmm3, xmm3, 8	vmovdqu xmm1, [rsi + 16]
	movzx	ecx, [3+rd1]	movzx ecx, [rd1+2]	vmovdqu xmm6, [rsi + 32]
	imui	ecx, eax	movzx eax, [rd1+3]	vmovaqu xmm/, [rsi + 48]
	add	rlid, [rdx]	vpackuswb xmm4, xmm10, xmm4	vmovdqa xmm2, [rip + .LCP10
	add	riid, r8d	vpackuswb xmm1, xmm1, xmm3	vpsnuib xmm3, xmm7, xmm2
	add	riid, riud	vpand xmm10, xmm0, xmm2	vpsnuib xmm2, xmm6, xmm2
	add	riid, ecx		vpunpektaq xnunz, xnunz, xnunz
	mov	[rax], riia	vpand xnuns, xnuno, xnuno	Villovada xillis, [rip + .LCP10
	movzx		vpsriw xninz, xninz, 8	vpshulb xnun4, xnun1, xnun3
	movsx	rlid, [4+rs1]	vmova xmm8, r9a	vpsnuib xmm3, xmm0, xmm3
	Imul	rya, rita	vpackuswb xnuns, xnunio, xnuns	vpunpektaq xnins, xnins, xnin4
	movsx	eax, [5+rs1]	vpsriw xmm6, xmm6, 8	vpblenda xmm12, xmm3, xmm2,
	mov2x		vpand xnunito, xnuno, xnun4	vinovada xinis, [rip + .LCPi0
	IMUI	rou, eax	vpackuswb xnunz, xnunz, xnuno	vpshulb xhun4, xhun7, xhun3
	movzx	eax, [2+rdi]	vpsriw xnun4, xnun4, 8	vpsnutb xnuns, xnuns, xnuns
	imul	ecx, [0+rsi]	vpand xnuno, xnuno, xnuno	vpunpekidq xmm3, xmm3, xmm4
	IMUI	eax, ecx	vpsriw xnuns, xnuns, s	Villovada xilili4, [rip + .LCPi0
	movsx	rid, [/tisi]	vpackuswo xillillo, xillillo, xillillo	wpshufb xmm4 ymm0 ymm4
	imul	riid, [Strai]	upackusub wmm6 wmm4 wmm3	wpuppekida wmm4 wmm4
	add	rod [4+rdy]	upand wmm3 wmm0 wmm1	wpblondd wmm3 wmm4 wmm3
	add	rod rod	upand xmm0, xmm0, xmm2	vpbiendo xmm4 [rip + ICPI0
el C Compiler	add	red oov	upackusub wmm3 wmm3 wmm0	woshufb wmm5 wmm7 wmm4
u C Compilor	add	rod r11d	uperida ymm0 ymm10 8	woshufb xmm4 xmm6 xmm4
	mou	[4+rdy] r9d	vpsiid xmm0, xmm10	wpuppekldg ymm4 ymm5
_VM compiler front-end	movzzy	r8d [rdi]	vpnovskow knuit, knuito	vmoudae vmm5 [rip + LCPI0
nsor Virtual Machine	movex	rod [8+rei]	womullw wmm4 wmm4	woshufb ymm2 ymm1 ymm5
	imul	r8d. r9d	vomovsybw xmm0, xmm0	vpshufb xmm5, xmm0, xmm5
lectorizer Generator	movsy	eav [9+rsi]	vpmovSxbw Xnuto, Xnuto	vounneklda ymm2 ymm5 ymm?
	DICTDIC		A: Approved for public release	distribution is unlimited
	UISIKIE		A. Approved for public release	wovdga xmm4, [rip + 1CP10

vmovdqa xmm5, [rip + .LCPI0 5]

Clang/LLVM (2.2x)

vmovdqu xmm0, [rsi] vmovdqu xmm1, [rsi + 16] vmovdqu xmm6, [rsi + 32] vmovdqu xmm7, [rsi + 48] vmovdga xmm2, [rip + .LCPI0 0] vpshufb xmm3, xmm7, xmm2 vpshufb xmm2, xmm6, xmm2 vpunpckldq xmm2, xmm2, xmm3 vmovdqa xmm3, [rip + .LCPI0 1]

vpshufb xmm4, xmm7, xmm3 vpshufb xmm3, xmm6, xmm3 vpunpckldg xmm3, xmm3, xmm4 vmovdqa xmm4, [rip + .LCPI0 3] vpshufb xmm5, xmm1, xmm4 vpshufb xmm4, xmm0, xmm4 vpunpckldq xmm4, xmm4, xmm5 vpblendd xmm3, xmm4, xmm3, 12 vmovdqa xmm4, [rip + .LCPI0 4]

vpunpckldq xmm3, xmm3, xmm4 vpblendd xmm12, xmm3, xmm2, 12 vmovdga xmm3, [rip + .LCPI0 2]

VeGen (11x)

vmovdqu64 zmm0, [rdx] vpbroadcastd zmm1, [rdi] vpdpbusd zmm0, zmm0, [rsi] vmovdqu64 [rdx], zmm0

Limitations: Only one type of optimization

DARPA Reason for confidence: Ithemal: Automatic construction of accurate processor performance models



Ithemal (Instruction THroughput Estimator using MAchine Learning) Mendis et al. [ICML'19]



Accuracy of Performance Model

Architecture	Predictor	Error		Method
Ivy Bridge	Default	33.5%		llvm-mca
	Diff Iune Ithemal	<u>25.4% ± 0.5%</u> 9.4%	F	IACA
	IACA OpenTuner	15.7% 102.0%		Ithemal
Haswell	Default	25.0% 23.7% ± 1.5%		Empirical execution
	Ithemal	9.2%		Bigger is
	IACA OpenTuner	17.1% 105.4%		
Skylake	Default DiffTune	26.7% 23.0% <u>+</u> 1.4%	Smaller is b	
	Ithemal	9.3%		
	IACA OpenTuner	14.3% 113.0%		

Throughput of Compiled Code

	Method	Throughput (Instructions / second)	
	llvm-mca	492	
	IACA	541	
	Ithemal	560	
	Empirical execution	13	
Bigger is better			
	Smaller is better	Limitations:Only one architectureOnly one figure of merit	

Renda, A., Chen, Y., Mendis, C., & Carbin, M. (2020, October). Difftune: Optimizing cpu simulator parameters with learned differentiable surrogates. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)* (pp. 442-455). IEEE.



Program Structure

- **1. ML-based Generators**: Multiple teams produce ML-based generators for different parts of compiler framework and one or two teams create a full ML-optimized tool chain
 - 1. Mapping
 - 2. Choosing and sequencing optimizations
 - 3. Performance Modeling appropriate for front end, middle, and back end
 - 4. Validation
 - 5. Code generation
 - 6. Integrated solutions: end-to-end compiler toolchains generated by the performer's and others' generators Full

Parts

- **2. Evaluation:** Government team creates a series of challenge problems of increasing complexity and measures component and system performance
 - 1. Using an ensemble of available hardware components and conventional languages
 - 2. Using models of emerging platforms
 - 3. Using hardware-agnostic languages
- **3. Consortium Development:** creates an institutional framework to incentivize industry adoption and sustainment of the technology



Program Timeline



- "Learning back-end performance models" (Back-End) doesn't rely on results in other areas and will provide capabilities for current systems (CPU's GPU's) in year 1 and will extend to other types of components in years 2 & 3
- "Learning to apply IR-based optimizations" (Middle) relies on back-end performance models and will provide capabilities in year 2 with extensions in year 3
- "Learning to allocate code segments to hardware components" (Front-End) depends on performance models from Middle. Final results in Year 3
- "End-to-end optimization" depends on all other technologies. Experimentation will proceed throughout with synthetic data. Final results in Year 3



Experimental design

- 1. Government evaluation support team creates a **"test suite"** of system software and test criteria (e.g., speed, power, code size)
- 2. Technical teams generate components for **Front-end**, **Middle**, and **Back-end** and integrate these into end-to-end solutions
- 3. Evaluation support team runs the **full spectrum systems** and **measures performance** and lines of compiler instructions needed to configure the full spectrum systems.





Notional Final Challenge Configuration



- The Year 3 challenge platforms will include an ensemble including 6 or more types of computational elements (and potentially more than 1 from each category)
- A variety of different configurations will be used as challenge platforms
- A variety of different programs will be compiled for these platforms



3 cycles of evaluation

- A mini test at 6-month boundaries (1.1, 2.1, 3.1)
- A full test at yearly boundaries (1.2, 2.2, 3.2)

Year 1 focuses on increasing the number of available accelerators

Year 2 add increasing complexity of the figure of merit (i.e., # of criteria – speed, power, code size)

Year 3 emphasizes the degree of hardware independence and application code size

	Eval #	# Component Types	Total # Components	Criteria in Figure of Merit	Hardware Agnostic	Source Code Size
	1.1	2	2	1	75%	10,000
	1.2	2	4	1	75%	25,000
	2.1	3	6	2	75%	50,000
	2.2	4	8	2	80%	100,000
5	3.1	5	10	3	85%	500,000
1	3.2	6	12	3	90%	1,000,000

MOCHA: ML and Optimization guided compilers that require minimal human effort to generate efficient code for heterogeneous hardware platforms

Year 3 Challenge problems include:

- Hardware agnostic languages
- All of the below

Year 2 Challenge problems include:

• Simulation models (or real systems) of emerging components in addition to those below

Year 1 Challenge problems include:

- Conventional programming languages
- Heterogeneous ensemble of existing hardware targets (CPUs, GPUs, ...)







- 1. Human effort reduction:
 - Time to create new compiler components
 - Measured by reduction in hand-written lines of compiler instructions
- 2. Code Quality multi-dimensional performance gain on benchmarks vs. baseline without ML guidance
 - Execution speed
 - Power consumption
 - Memory footprint

Metric	Year 1	Year 2	Year 3
Human Effort Reduction	50%	75%	90%
Code quality performance gain	1x	2.5x	5x



www.darpa.mil