



# Static Performance Analysis with LLVM

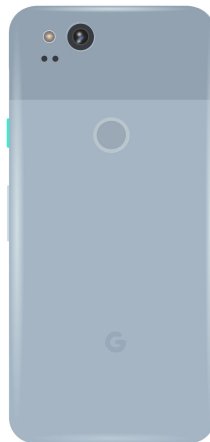
**Clément Courbet**

G. Chatelet, B. De Backer, O. Sykora,

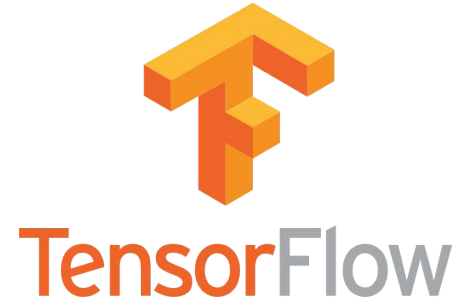
Google Compiler Research



TensorFlow



Low-precision matrix multiplication  
[github.com/google/gemmlowp](https://github.com/google/gemmlowp)



**Optimize this at (nearly) all costs !**

# Tools

## Benchmarks

- Closer to real-life performance
- Sloooooow
- Requires access to hardware

## Static Analysis

- Fast
- Reproducible
- Hard to model input-dependent behaviour (branches, cache)

# Our Static Performance Analyzer

## MC Basic Block

```
1:  
movdqu (%rdi),%xmm1  
lea 0x10(%rdi),%rdi  
movdqu (%rsi),%xmm2  
lea 0x10(%rsi),%rsi  
movdqa %xmm1,%xmm3  
psubusb %xmm2,%xmm1  
psubusb %xmm3,%xmm2  
por %xmm2,%xmm1  
movdqa %xmm1,%xmm2  
punpcklbw %xmm5,%xmm1  
punpckhbw %xmm5,%xmm2  
pmaddwd %xmm1,%xmm1  
pmaddwd %xmm2,%xmm2  
padd %xmm1,%xmm0  
padd %xmm2,%xmm0  
sub $0x10,%ecx  
jg 1b
```

Simulator

Annotated Trace

Port Pressure

Latency/ Inverse Throughput



Scheduler

🤖 No semantics 🤖

# Simulator API

Target-independent [Simulator](#) interface:

```
const vector<MCInst>& BasicBlock = ...;  
auto Simulator = Target->createSimulator();  
SimulationLog Log = Simulator.Run(BasicBlock);
```

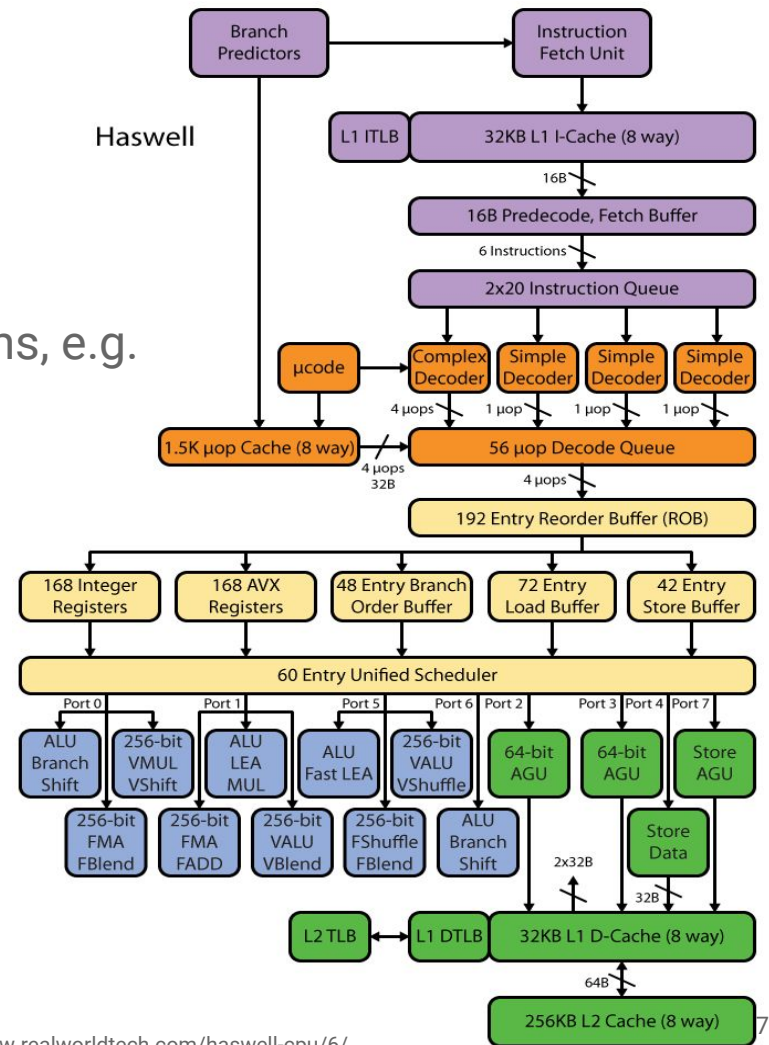
# Simulator Internals: Components

- Generic →  
Reuse LLVM's target-independent descriptions, e.g.

[Scheduler](#): `llvm::MCSchedModel`

[RegisterRenamer](#): `llvm::MCRegisterInfo`

- Target-specific, e.g. Intel [Fetcher](#)



analyzing 'libyuv\_sumsquareerrorsse2.s'  
 ran 20 iterations in 132 cycles  
 Block Inverse Throughput: [5-6] cycles per iteration

# Analysis: IACA-like frontend

Port Pressure (cycles per iteration):

Port	HWDivider	HWP0	HWP1	HWP2	HWP3	HWP4	HWP5	HWP6	HWP7
Cycles		4.35	4.40	1.00	1.00		4.30	1.95	

#Ops	HWDivider	HWP0	HWP1	HWP2	HWP3	HWP4	HWP5	HWP6	HWP7
1				1.00					
1			0.30				0.70		
1					1.00				
1			0.55				0.45		
1		0.35	0.65						
1			0.70				0.30		
1			0.40				0.60		
1		0.95	0.05						
1		1.00							
1							1.00		
1							1.00		
1		1.00							
1		1.00							
1			0.75				0.25		
1			1.00						
1		0.05						0.95	
1								1.00	

```

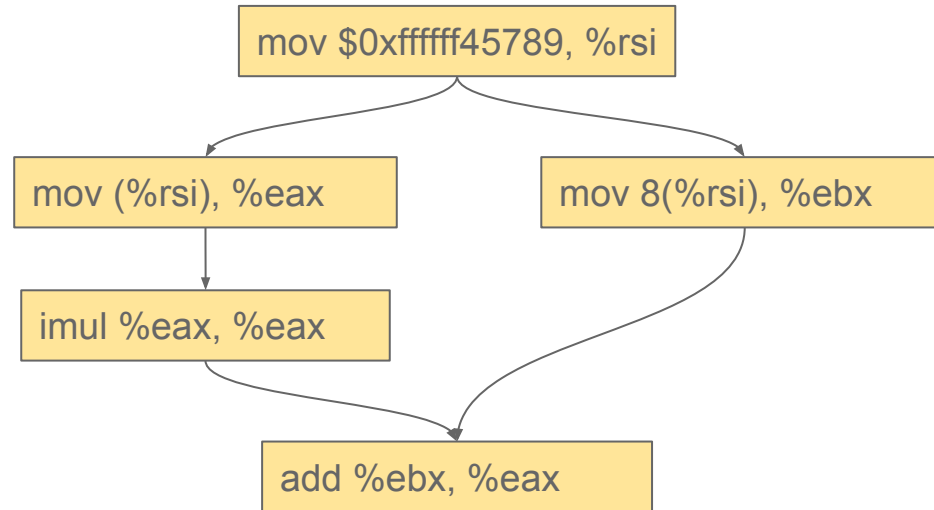
movdqu xmm1, xmmword ptr [rdi]
lea rdi, [rdi + 0x10]
movdqu xmm2, xmmword ptr [rsi]
lea rsi, [rsi + 0x10]
movdqa xmm3, xmm1
psubusb xmm1, xmm2
psubusb xmm2, xmm3
por xmm1, xmm2
movdqa xmm2, xmm1
punpcklbw xmm1, xmm5
punpckhbw xmm2, xmm5
pmaddwd xmm1, xmm1
pmaddwd xmm2, xmm2
padd xmm0, xmm1
padd xmm0, xmm2
sub ecx, 0x10
jg .Ltmp0
  
```



# Automatic Scheduling

Minimize the simulated latency (alternative to PostRAMachineSched)

- Random/Exhaustive Search ([CP Solver](#)) ~hours-days
- Genetic Algorithms (Biased Random-Key Genetic Algorithms) <seconds



# Exhaustive Search

- gemmlowp's SSE4\_32\_Kerne14x4Depth2:  
0-2% faster (vs. implementation contributed by Intel).

- libwebp's FTransform():  
0-5% faster.

On benchmarks; no performance regressions



# Genetic Algorithms

- 100 milliseconds
- 10% improvement (in theory)

Original: **9.5** Cycles/Iter

```
movd mm0, dword ptr [r8]
punpcklbw mm0, mm7
movq mm1, mm0
movq mm2, mm0
pmullw mm1, mm1
paddw mm1, mm3
psrlw mm1, 0x8
pmullw mm0, mm1
paddw mm0, mm3
psrlw mm0, 0x8
pmullw mm0, mm4
pmullw mm1, mm5
pmullw mm2, mm6
paddw mm0, mm1
paddw mm0, mm2
psrlw mm0, 0x6
packuswb mm0, mm7
movd dword ptr [r8], mm0
add r8, 0x4
dec ecx
jne .Ltmp0
```

Rescheduled: **8.5** Cycles/Iter

```
movd mm0, dword ptr [r8]
punpcklbw mm0, mm7
movq mm1, mm0
pmullw mm1, mm1
movq mm2, mm0
pmullw mm2, mm6
paddw mm1, mm3
psrlw mm1, 0x8
pmullw mm0, mm1
pmullw mm1, mm5
paddw mm0, mm3
psrlw mm0, 0x8
pmullw mm0, mm4
paddw mm0, mm1
paddw mm0, mm2
psrlw mm0, 0x6
packuswb mm0, mm7
movd dword ptr [r8], mm0
add r8, 0x4
dec ecx
jne .Ltmp0
```

# Future Work

- Integrating into [llvm-mca](#) (in particular frontend simulation)
- Genetic scheduler → MachineFunctionPass

# Try It Out!

[https://github.com/google/EXEgesis/tree/master/llvm\\_sim](https://github.com/google/EXEgesis/tree/master/llvm_sim)