



### **Table of Contents**

- Moore's Law Novel Hardware Architectures
- Performance Engineering
- DOCC A self-learning Compiler
- Daisytuner Optimizing Compiler Collection
  - Tuning
  - Auto-Tuning
  - Transfer Tuning
- Outlook



#### Moore's Law: "The number of transistors doubles every two years"

1970-2000:2x Speedup every twoyears on the same code

2000+: What's happening now?

CPU	Year	Transistors	Clock	Structure	Caches
4004	1971	2,300	740 kHZ	20 Micro	
8008	1972	3,500	500 kHz	10 Micro	
8086	1978	29,000	10 MHz	3 Micro	
80286	1982	134,000	25 MHz	1.5 Micro	
80386	1985	275,000	33 MHz	1 Micro	
80486	1989	1,2M	50 MHz	0.8 Micro	8К
Pentium I	1994	3.1M	66 MHz	0.8 Micro	8К
Pentium II	1997	7.5M	300 MHz	0.35 Micro	16K/512K
Pentium III	1999	9.5M	600 MHz	0.25 Micro	16K/512K
Pentium IV	2000	42M	1.5 GHz	0.18 Micro	8K/256K
P IV F	2005	~160M	2.8 GHz	90 nm	16K/2MB
Core I7	2008	781M	3.2 GHz	45 nm	32/256/8
Ivy Bridge	2013	2,890M	3.3 GHz	22 nm	32/256/30

3



Problem:

Power = clock<sup>3</sup>

-> Top 500 consume 3 Billion
kWh per year
-> Use of AI & digital twins
leads to increasing resource
needs



Solution:

Const: frequency Grow: cores (and caches)



#### Adapt your code to the hardware!



Solution:



Const: frequency Grow: cores (and caches) Specialization: accelerated computing





Adapt your code to the hardware!



Solution:



Const: frequency Grow: cores (and caches) Specialization: accelerated computing





Performance is a primary design goal. Rewriting code is expensive!









Optimized HPC Application



















A long term consideration would be to introduce a different memory layout, with better cache locality and clearer access patterns for the compiler to vectorize with more ease. Figure 6 depicts the current memory layout and a reordered



A long term consideration would be to introduce a different memory layout, with better cache locality and clearer access patterns for the compiler to vectorize with more ease. Figure 6 depicts the current memory layout and a reordered

Can we talk to the compiler (to generate those insights automatically) ?



# DOCC: A compiler that learns with every program (based on LLVM plugins, cf. transfer tuning)







#### **Continuous Benchmarking**

Partition	CPU	Accelerators	Description
bellis4	ARM Cortex A-72	Coral Edge TPU (USB)	Raspberry Pi 4
bellis5	ARM Cortex A-76	HAILO 8L	Raspberry Pi 5
tansy	ARM Cortex A-78	NVIDIA Ampere	NVIDIA Jetson Orin
			Nano Super
zinnia	AMD EPYC 9554	NVIDIA A16, Tenstorrent Wormhole	
		N150, Axelera Al Metis	
chamomile	Intel Xeon Silver	NVIDIA A10, Tenstorrent Wormhole	
	4516Y+	N150, Axelera Al Metis	
aster	Ampere Altra	NVIDIA L4, Tenstorrent Wormhole	
(soon)	Max M128-30	N150, Axelera Al Metis	

NOTE: We're currently working on bringing up the Axelera AI Metis and HAILO 8L processors.



daisytuner.com





00 	Benchmar	ks	00	
Benchmark	Runtime	Diff	Threshold	
# matmul	7.11 ms	-3.67%	N/A	
# matmul2	7.58 ms	+0.64%	N/A	



#### DOCC

# DOCC: A compiler that learns with every program (based on LLVM plugins, cf. transfer tuning)





### **Table of Contents**

- Moore's Law Novel Hardware Architectures
- Performance Engineering
- DOCC A self-learning Compiler
- Daisytuner Optimizing Compiler Collection
  - Tuning
  - Auto-Tuning
  - Transfer Tuning
- Outlook





#### Tuning

<u>ICS '23: Proceedings of the 37th International Conference on Supercomputing</u> • Pages 50 - 62 <u>https://doi.org/10.1145/3577193.3593714</u>







#### Auto-Tuning



Challenge:

- Sensible optimization schema
- "Good" optimization function
- Ensurance of correctness



#### Auto-Tuning



Challenge:

- Sensible optimization schema
- "Good" optimization function
- Ensurance of correctness

#### **Initial state**





#### Auto-Tuning



Challenge:

- Sensible optimization schema
- "Good" optimization function
- Ensurance of correctness

#### Initial state

#### **Target state**

























Vector Key \_ (embedding)

#pragma omp parallel for for (int i=0; i < 1024; i++) for (int j=0; j < 1024; j++) for (int k=0; k < 1024; k++) C[i][j] += A[i][k] \* B[k][j]

Loop Nest



<sup>1</sup> Ben-Nun et. al. - Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures





<sup>1</sup> Ben-Nun et. al. - Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures









Performance Portability on Heterogeneous Architectures







- Generate new code
- Introduce source-level instrumentation (PAPI)
- Measure representative performance counter

Processo	r
Branch Misprediction Ratio:	0.02
Cycles Per Instruction:	0.44
FLOP [DP]:	8.28 MFLOP
Runtime:	19.10 ms
+	+
1	
+	+
Caches	
+	+
L2 Volume:	113.16 MB
+	+
1	
+	••••••
Main Mem	ory
+	
Memory Volume:	32.36 MB
+	+





36





37















#### **Performance Model**

Challenge:

- Number of optimizable regions in practice high
- Particularly in engineering, differing program flow depending on input
- Where is the biggest optimization potential?



#### **Performance Model**

Challenge:

- Number of optimizable regions in practice high
- Particularly in engineering, differing program flow depending on input
- Where is the biggest optimization potential?

- Find bottlenecks
- Approximate optimization potential



### **Performance Model**



daisytuner (bot) commented 2 weeks ago

#### Daisytuner Report - polybench (raspi5)

Hotspots



Regions	Operational Intensity	Peak Performance	Relative Runtime
Region starting at line 112 in tests/polybench/datamining/correlation/correlation.c	0.39 FLOP/Byte	19.57%	47.84%
Region starting at line 103 in tests/polybench/datamining/correlation/correlation.c	0.50 FLOP/Byte	33.99%	0.69%
Region starting at line 90 in tests/polybench/datamining/correlation/correlation.c	0.18 FLOP/Byte	68.37%	0.49%
Region starting at line 81 in tests/polybench/datamining/correlation/correlation.c	0.08 FLOP/Byte	74.41%	0.51%
Region starting at line 36 in tests/polybench/datamining/correlation/correlation.c	24.00 FLOP/Byte	18.47%	2.60%

Benchmark ()

#### Measured Files







# What is Loop Scheduling

# 

HPC Application

Focus on Loop

Nests

```
for (int i=1; i < 10000; i++)
        A[i] = B[i-1] + B[i] + B[i+1];</pre>
```

# What is Loop Scheduling



46

# What is Loop Scheduling





## Normalization Criteria - Intuition

#### Matmul B (k swap)

```
for (i = 0; i < N; i++) {
  for (k = 0; k < K; k++) {
    for (j = 0; j < M; j++) {
        C[i][j] += A[i][k] * B[k][j];
    }
}</pre>
```

#### Matmul C (with init)

48

```
for (i = 0; i < N; i++) {
   for (j = 0; j < M; j++) {
     double sum = 0;
     for (k = 0; k < K; k++) {
        sum += A[i][k] * B[k][j];
     }
     C[i][j] = sum;
}</pre>
```

# A Priori Loop Nest Normalization - Idea

```
Algorithmic
```

```
for (int i = 0; i < N; i++) {
    for (int j = 0; j < M; j++) {
        double sum = 0;
        for (int k = 0; k < K; k++) {
            sum += A[i][k] * B[k][j];
        }
        C[i][j] = sum;
    }
</pre>
```

Language-specific Variation

C = A @ B

#### **Pre-Optimized**

# Matching Program Transformations





# Performance Embeddings - Evaluation



Model	Main Memory Bandwidth	Data Locality
Reuse Distance	0.78	0.87
IR2Vec <sup>1</sup>	0.47	0.41
Tiramisu <sup>2</sup>	0.32	0.35
Our Model	0.25	0.31

Mean variation of neighbors for different performance metrics. A lower value is better.

51

<sup>1</sup> VenkataKeerthy et al. - IR2Vec: LLVM IR based Scalable Program Embeddings

<sup>2</sup> Baghdadi et al. - A Deep Learning Based Cost Model for Automatic Code Optimization



#### Evaluation

Benchmark	Explored States	Reference [ms]	Transfer Tuning (k=5)	Transfer Tuning (k=10)
mlp	111,508	1.47	+ 38.8%	+ 37.4%
softmax	183,427	110.40	+ 0.6%	+ 0.5%
blur filter	1,342	1.03	0.0	0.0
daubechies wavelet	9,101	8.73	- 3.7%	- 92.0%
haar wavelet	8,639	0.22	0.0%	0.0%
harris filter	1,651	9.06	+ 0.2%	- 4.0%
histogram filter	147,438	32.51	+ 1.2%	- 4.9%
unsharpening filter	25,080	29.66	+ 3.1%	+ 0.5%
heat 3D	69,080	13428.98	+ 3.6%	+ 2.8%
horizontal diffusion	34,534	7.00	+ 4.8%	+ 2.8%
matmul	65,986	14.17	+ 5.3%	+ 4.1%
min-plus mm	65,999	24.76	+ 9.0%	+ 8.5%

Number of explored states and optimized runtime of the Auto-Scheduler compared to Transfer Tuning a constant number of hypotheses.



#### Transfer Tuning: CLOUDSC



(a) Strong scaling behavior of CLOUDSC for the Fortran, C, DaCe, and daisy versions from left to right and grouped by the number of threads.