The background of the slide features a large, semi-transparent watermark of the University of Delaware seal. The seal is circular and contains the text 'UNIVERSITY OF DELAWARE' around the perimeter, '1743' at the bottom, and 'SOL' in the center. It also depicts an open book with the words 'GRAMM', 'METAPH', 'PHILOL', 'LOGICA', 'RHETORICA', and 'ETHICA' on its pages.

Using Graph-Based Characterization for Predictive Modeling of Vectorizable Loop Nests

William Killian

PhD Preliminary Exam Presentation

Department of Computer and Information Science

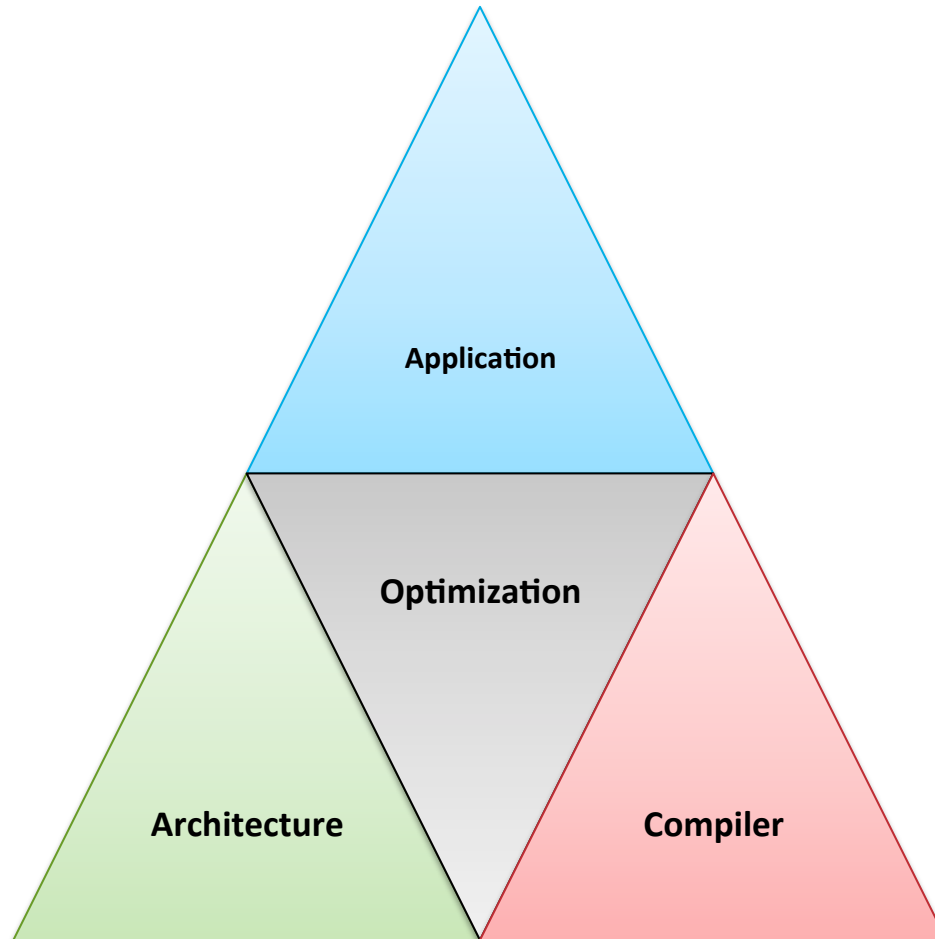
Committee

John Cavazos and Xiaoming Li

January 20, 2015



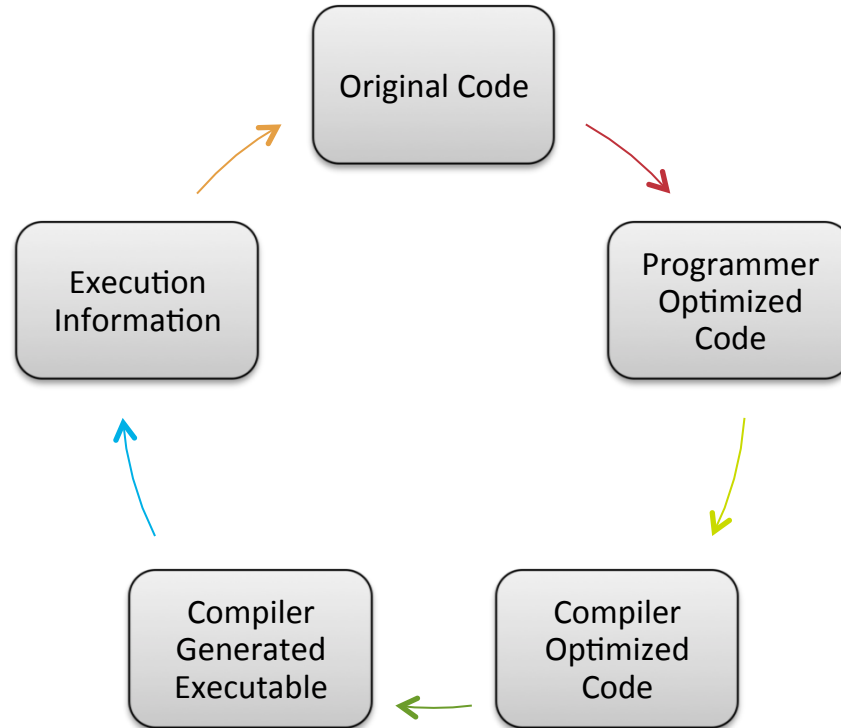
Code Optimization



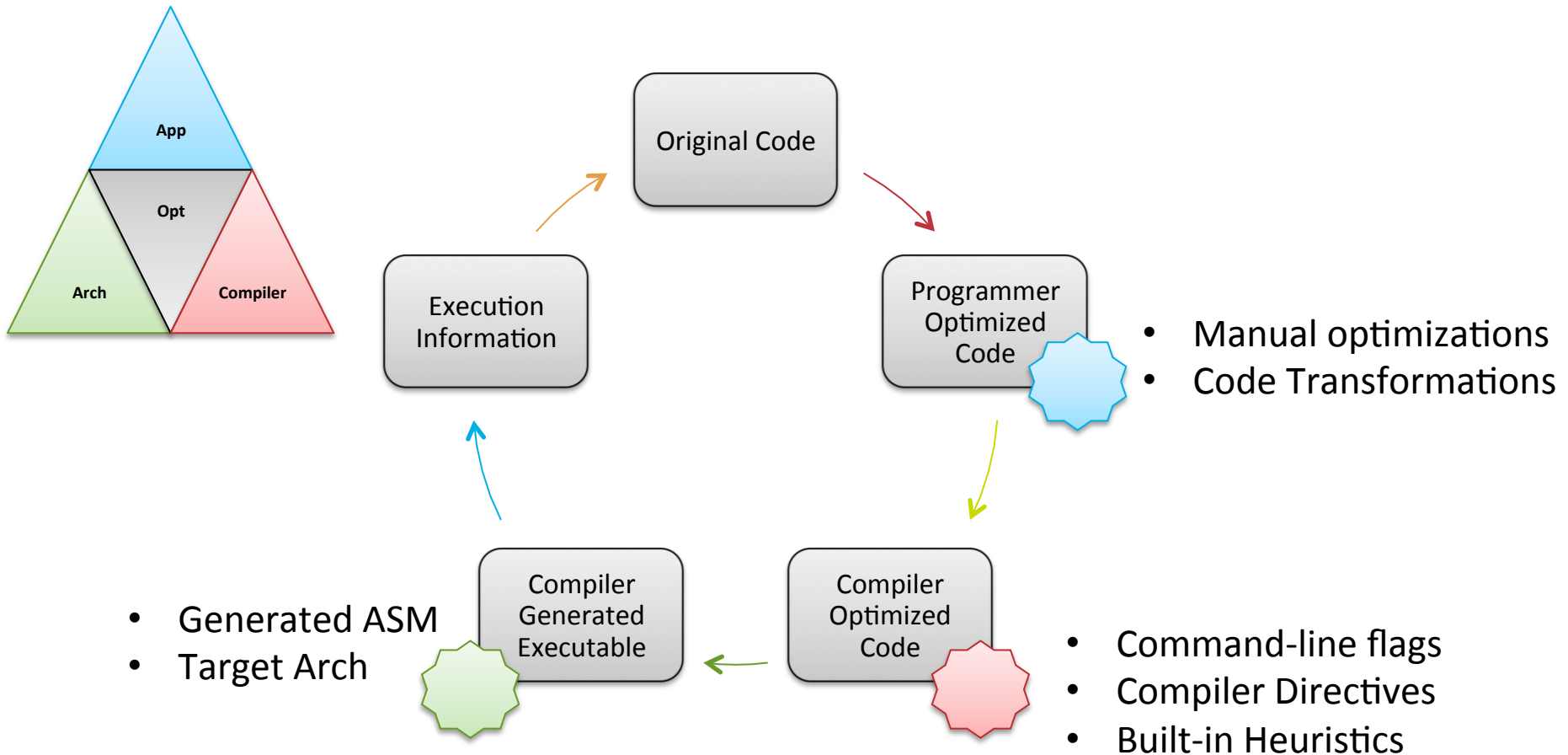
Problems with Optimizing Code

- New compilers
 - New optimizations
 - Extended language features
- New architectures
 - Old programming model won't work well (GPUs)
 - New/improved capabilities (ISA, cache coherency)
 - Compilers don't update with the architecture
- Old code
 - Legacy code expected to work
 - Maintenance of existing code

Workflow



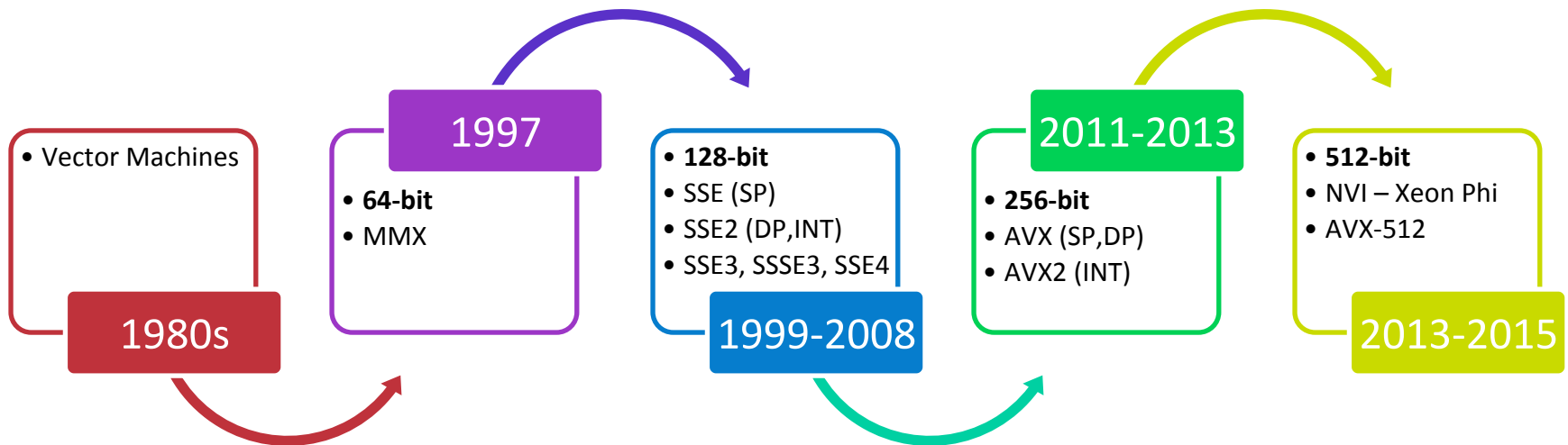
Workflow with Optimizations



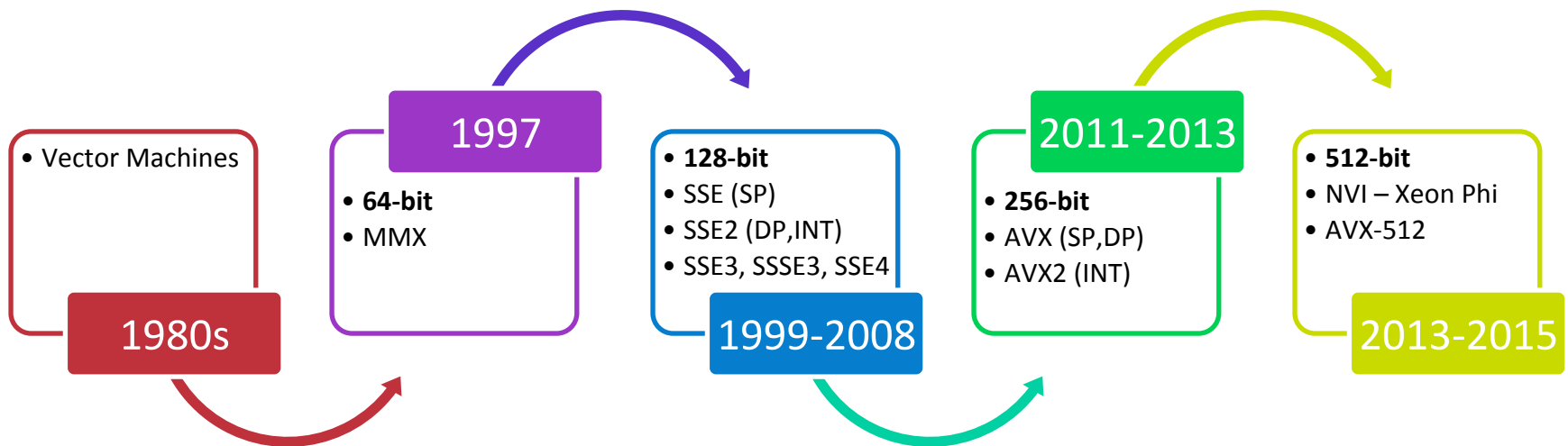
Optimizations

- Manual
 - Loop transformations – unrolling, fusion/fission
 - Data structure changes – Array of Structures → Structure of Arrays
- Compiler Directives
 - Source code hints to the compiler indicated by user
 - Usually used for local (scope or loop) optimizations
 - May automatically transform code (e.g. `#pragma unroll 4`)
- Command-Line flags
 - Optimizations applied to entire program (`-O3`, `-funroll=4`)
 - May specify target architecture and features permitted (e.g. `-march=avx`)

(Intel) SIMD Architecture Evolution



(Intel) SIMD Architecture Evolution



How can we choose the best optimizations to exploit vectorization?

Research Questions

Optimization Search Space:

With many types of vectorization optimizations, how do we choose which ones to apply?

Automation:

How can we automatically select optimizations, apply optimizations, and evaluate performance?

Automatic Performance Improvement:

How can we quickly select good vectorization optimizations that improve performance?

Research Questions

Optimization Search Space:

With many types of vectorization optimizations, how do we choose which ones to apply?

Automation:

How can we automatically select optimizations, apply optimizations, and evaluate performance?

Automatic Performance Improvement:

How can we quickly select good vectorization optimizations that improve performance?

Optimization Search Space

- Used source code directives to drive the optimization selection and modification
- Guide the internal compiler vectorization heuristics to improve performance.
- Use 6 different optimizations that provide varying levels of guidance to the compiler
- Exhaustive search space of directives on individual loop nests

Optimization Search Space

- Apply No Optimization
 - Let the compiler perform default vectorization
- `#pragma vector always` - May generate slower code
 - Ignore speed-up factor predicted by internal model
- `#pragma ivdep` - May generate invalid code
 - Ignore built-in check for all **unproven** vector dependences
 - Proven vector dependences will not be vectorized
- `#pragma simd` - May generate invalid code
 - Ignore all dependencies and reductions
 - Can vectorize an entire loop nest (outer-loop vectorization)
 - Optional argument `vectorlength(n)`. Vector length states how many safe iterations can be done at once ($n = 2, 4, 8$)

Research Questions

Optimization Search Space:

With many types of vectorization optimizations, how do we choose which ones to apply?

Automation:

How can we automatically select optimizations, apply optimizations, and evaluate performance?

Automatic Performance Improvement:

How can we quickly select good vectorization optimizations that improve performance?

Version Generation Automation

- Creation of two utilities
- autovec
 - Simplified directive language; provides support for permutation of optimizations
 - Source-to-source compiler
- VALT – Vectorization And Loop Transformation
 - Provides developer with concise language to specify vectorization and loop optimization directives
 - Extension of autovec
 - Supports multiple backend compilers

autovec language to Intel directives

autovec directive	Intel-specific pragma
<code>#pragma autovec permute</code>	Generates each of the following version into new file
<code>#pragma autovec vl(x)</code>	<code>#pragma simd vectorlength(x)</code>
<code>#pragma autovec ivdep</code>	<code>#pragma ivdep</code>
<code>#pragma autovec always</code>	<code>#pragma vector always</code>
<code>#pragma autovec none</code>	

Permute was configured to generate:

- No optimization
- `#pragma vector always`
- `#pragma ivdep`
- `#pragma simd vectorlength(2)`
- `#pragma simd vectorlength(4)`
- `#pragma simd vectorlength(8)`

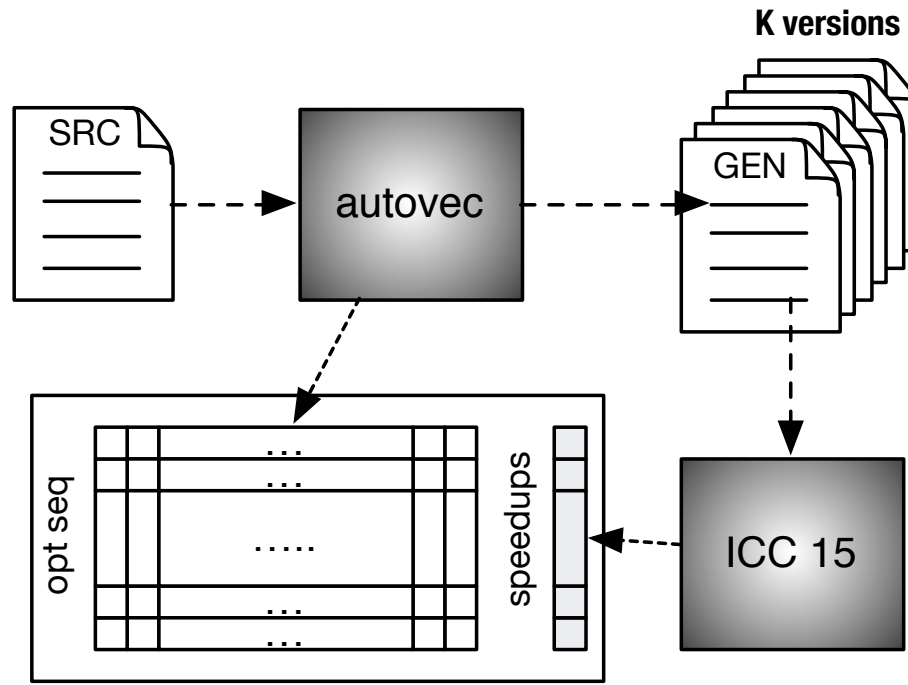
VALT language grammar

<p> $\langle \text{directive} \rangle ::= \text{\#pragma} \text{\ 'VALT' } \langle \text{clauselist} \rangle$ $\langle \text{clauselist} \rangle ::= \langle \text{clauselist} \rangle [[,]] \langle \text{clause} \rangle$ $\langle \text{empty} \rangle$ $\langle \text{clause} \rangle ::= \text{\ 'vector' } \text{\ (' } \langle \text{vectorlist} \rangle \text{\ ')}$ $\text{\ 'depend' } \text{\ (' } \langle \text{dependopts} \rangle \text{\ ')}$ $\text{\ 'vectorsize' } \text{\ (' } \langle \text{number} \rangle \text{\ ')}$ $\text{\ 'loop' } \text{\ (' } \langle \text{looplist} \rangle \text{\ ')}$ $\langle \text{vectorlist} \rangle ::= \text{\ 'none'}$ $\langle \text{vectorlist} \rangle \text{\ ', ' } \langle \text{vectoritem} \rangle$ $\langle \text{vectoritem} \rangle$ $\langle \text{looplist} \rangle ::= \langle \text{looplist} \rangle \text{\ ', ' } \langle \text{loopitem} \rangle$ $\langle \text{loopitem} \rangle$ </p>	<p> $\langle \text{vectoritem} \rangle ::= \text{\ 'default'}$ \ 'none' \ 'always' \ 'aligned' \ 'unaligned' \ 'nontemp' \ 'temp' $\langle \text{loopitem} \rangle ::= \text{\ 'unroll' } [[\text{\ (' } \langle \text{number} \rangle \text{\ ') }]]$ $\text{\ 'jam' } [[\text{\ (' } \langle \text{number} \rangle \text{\ ') }]]$ \ 'dist' \ 'nofusion' $\langle \text{dependopts} \rangle ::= \text{\ 'ignore'}$ \ 'default' $\langle \text{number} \rangle ::= [1-9][0-9]^*$ </p>
--	--

VALT language to Intel directives

VALT directive	Intel-specific pragma
<code>#pragma vector(default)</code>	No code emitted
<code>#pragma vector(none)</code>	<code>#pragma novector</code>
<code>#pragma vector(always)</code>	<code>#pragma vector always</code>
<code>#pragma vector(ignore)</code>	<code>#pragma ivdep</code>
<code>#pragma vector(aligned)</code>	<code>#pragma vector aligned</code>
<code>#pragma vector(temp)</code>	<code>#pragma vector temporal</code>
<code>#pragma vector(nontemp)</code>	<code>#pragma vector nontemporal</code>
<code>#pragma vectorsize(x)</code>	<code>#pragma simd vectorlength(x)</code>
<code>#pragma loop(unroll(x))</code>	<code>#pragma unroll(x)</code>
<code>#pragma loop(jam(x))</code>	<code>#pragma unroll_and_jam(x)</code>
<code>#pragma loop(nofusion)</code>	<code>#pragma nofusion</code>
<code>#pragma loop(dist)</code>	<code>#pragma distribute_point</code>

Version Generation Workflow



Research Questions

Optimization Search Space:

With many types of vectorization optimizations, how do we choose which ones to apply?

Automation:

How can we automatically select optimizations, apply optimizations, and evaluate performance?

Automatic Performance Improvement:

How can we quickly select good vectorization optimizations that improve performance?

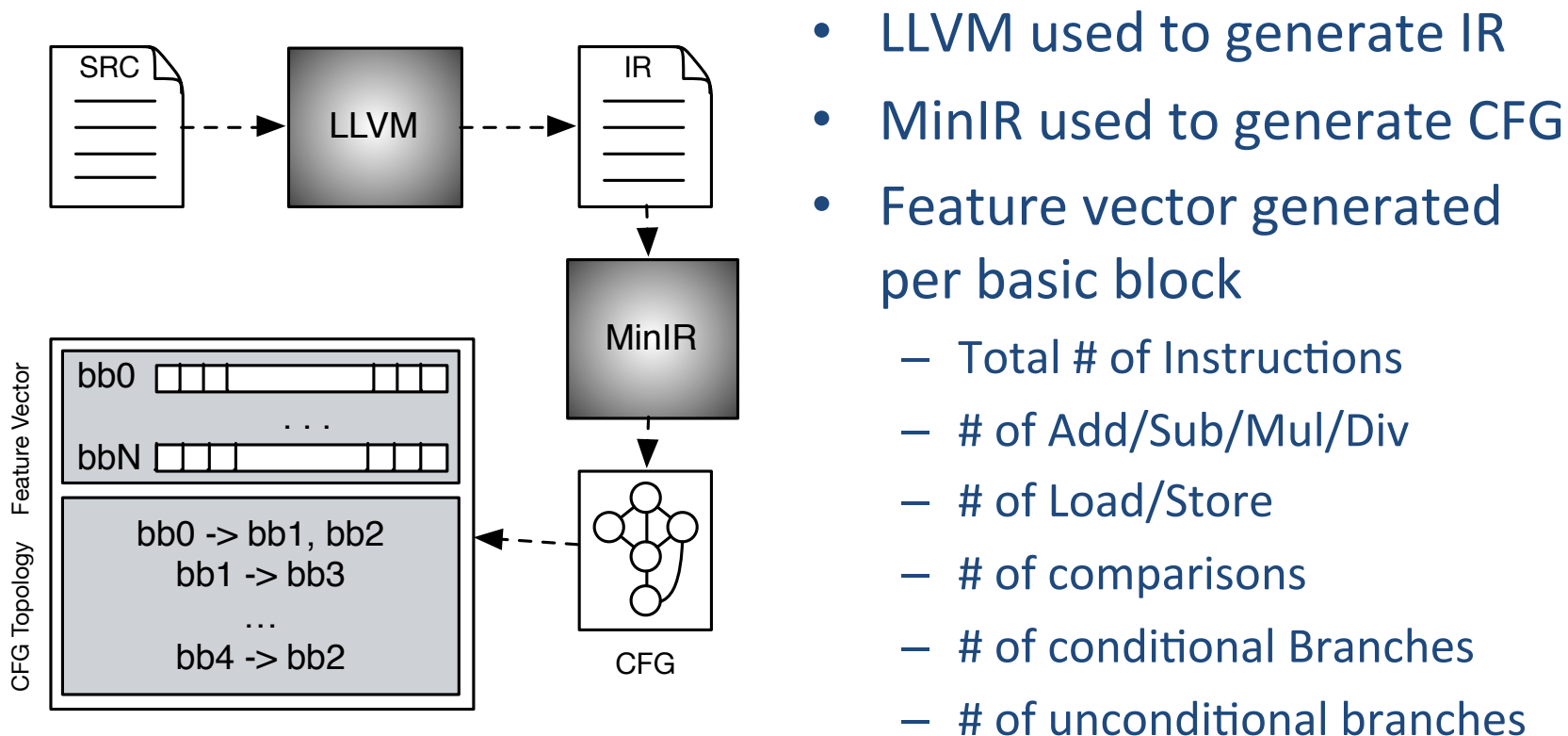
Machine Learning - Previous Solutions

- Stock et al. proposed using machine learning techniques to improve automatic vectorization
- Park et al. proposed using graph-based learning techniques to optimize programs at loop-nest granularity

Proposed Solution:

- Use graph-based learning techniques to choose vectorization optimizations for vectorizable loop nests
- Construct a graph-based speedup predictor that can predict a speedup when applying vectorization optimizations to a loop nest

Feature Extraction



Example Control Flow Graph for Loop Nest

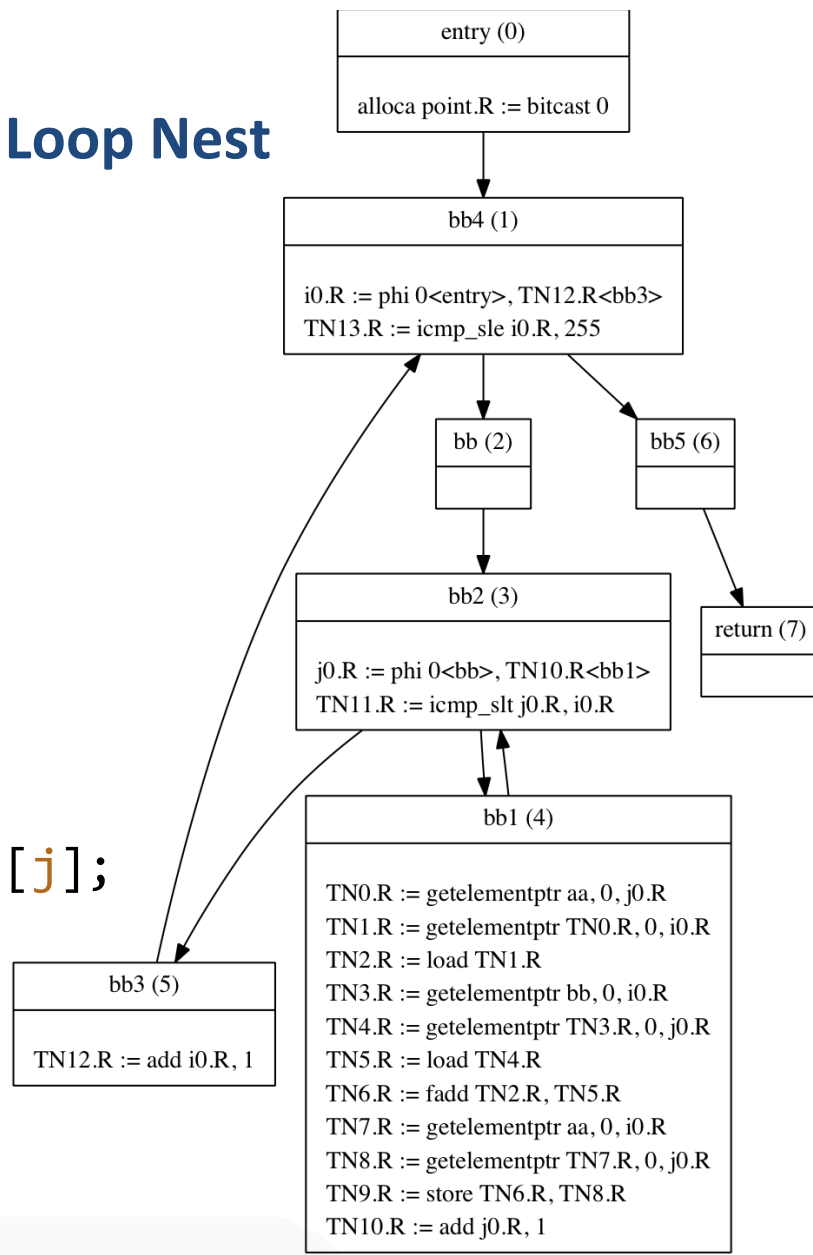
- Total of 6 basic blocks
- One entry, one return
- Basic blocks may not contain code

```
float aa[LEN2][LEN2];
```

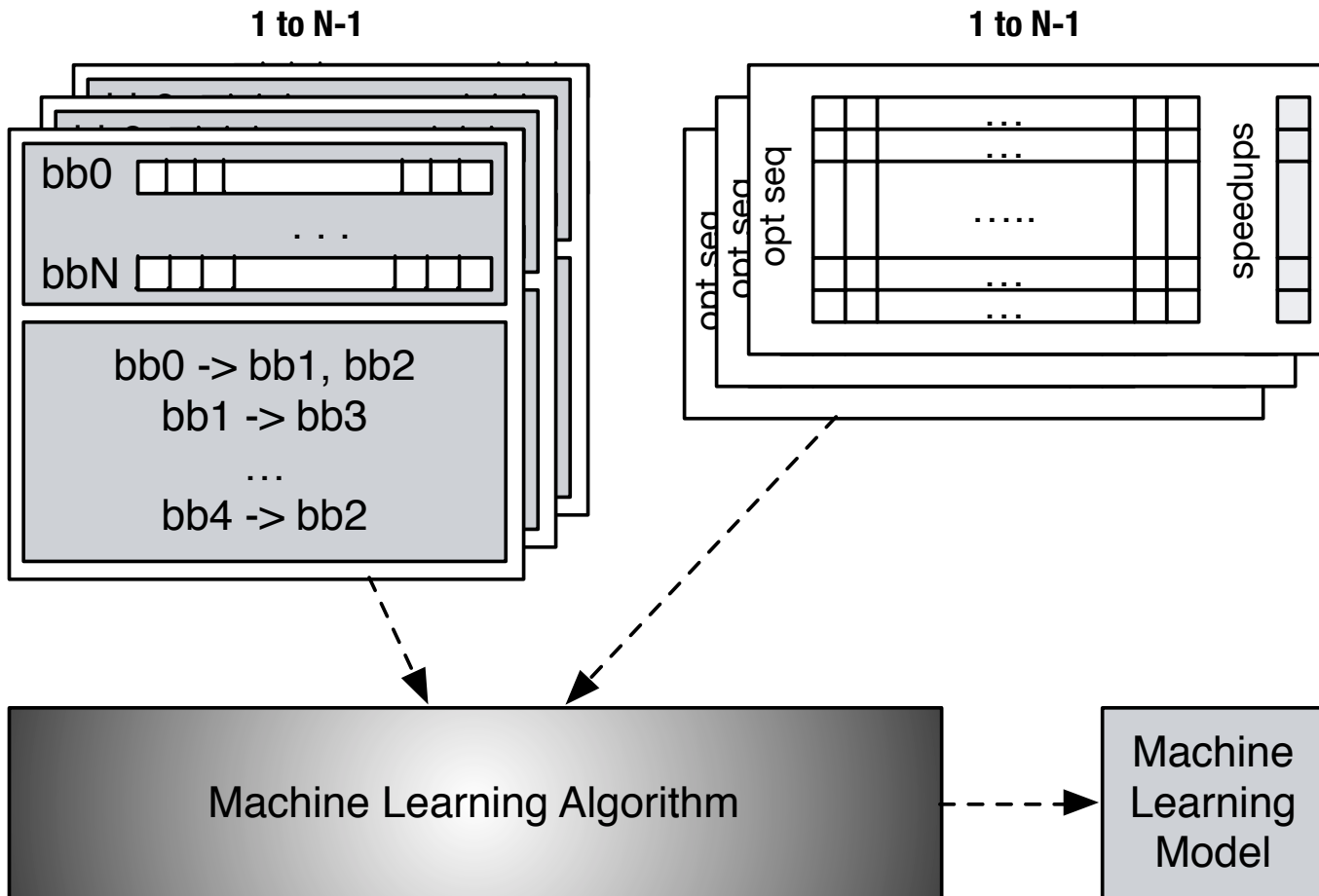
```
for (int i = 0; i < LEN2; i++)
```

```
  for (int j = 0; j < i; j++)
```

```
    aa[i][j] = aa[j][i] + bb[i][j];
```



Machine Learning Model Construction



Optimization Encoding

Bit Configuration	Encoded Optimization
000000	No Loop
100000	No Optimization Performed
110000	<code>#pragma vector always</code>
111000	<code>#pragma ivdep</code>
111100	<code>#pragma simd vectorlength(2)</code>
111110	<code>#pragma simd vectorlength(4)</code>
111111	<code>#pragma simd vectorlength(8)</code>

Machine Learning Algorithm

1. Training data:

L = set of all loop nests

O = set of all optimization sequences

$\text{speedup}(l, o)$ - observed speedup from applying o to loop nest l

$\text{scores} = \{(l, o, \text{speedup}(l, o)) \mid \forall l \in L, \forall o \in O, l_{\text{size}} = o_{\text{size}} \wedge \text{valid}(l, o)\}$

2. Construct kernel similarity matrix by computing similarity between all training data points (loop nest + optimization)

$\forall i \text{ in } [0, \|\text{scores}\|), \forall j \text{ in } [0, \|\text{scores}\|)$

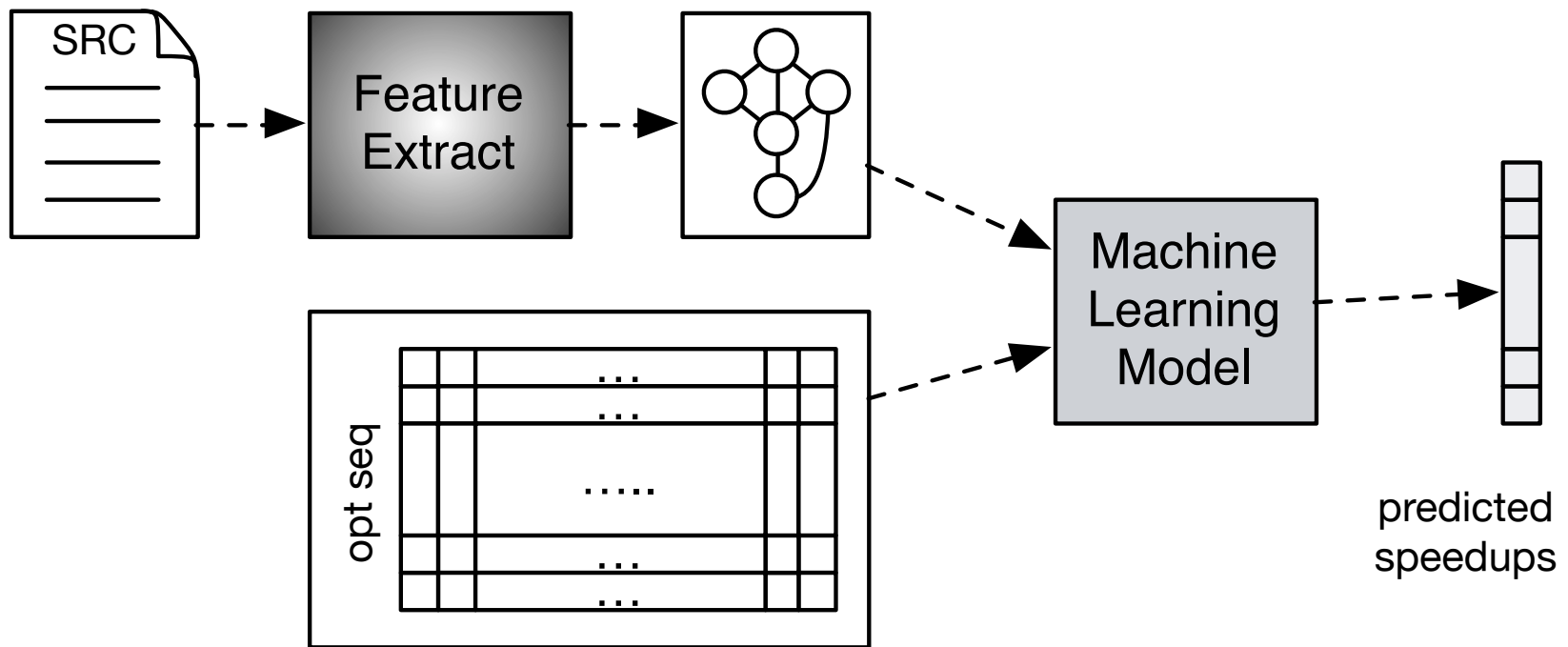
$$\text{Sim}_{km}^{i,j} = \text{Sim}_{loop}^{i,j} \times \|\text{count1}(\text{scores}_o^i) - \text{count1}(\text{scores}_o^j)\|$$

3. Train on kernel matrix with speedup as production target

Machine Learning Algorithm

- *Graph kernel* functions used to transform the training into a different, linearly-separable feature space.
 - Shortest path graph kernel
 - Used previously by Park et al. for speedup predictors
 - Similarity calculated by normalizing intersection kernel matrix
- Linear classifier is constructed that separates the points into multiple classes.
- Support vector machines (SVMs) used to construct predictive models from the kernel similarity matrix
- Predicts a speedup

Using Machine Learning Model for Unseen Program



EXPERIMENT SETUP

Benchmark Selection

TSVC

- 151 loop nests with varying access patterns, computations, and memory access types
- Originally used to evaluate how well a compiler can recognize patterns for vectorization
- Millisecond timing granularity with each loop-nest within a repeat loop

PolyBench/C

- 30 static control-flow micro-benchmarks from several scientific domains
- Modified to create different individual versions optimizing single loop nests (65 total loop nests)
- Clock-tick timing granularity across the entire kernel execution

Machine Configuration

Nehalem (NHM)

- Core i7 950
- 3.06GHz quad-core
- 8MB L3 cache
- 24GB DDR3-1333
- 128-bit vector width
- Up to SSE 4.2 ISA
- 45nm
- Q2 2009

Haswell (HSW)

- Core i7 5930K
- 3.5GHz hex-core
- 15MB L3 cache
- 32GB DDR4-2133
- 256-bit vector width
- Up to AVX2 ISA
- 22nm
- Q4 2013

Processor dynamic frequency scaling was disabled for all experiments
We only analyzed speedup for cross-architecture comparison, not performance

Machine Configuration

Nehalem (NHM)

- Core i7 950
- 3.06GHz quad-core
- 8MB L3 cache
- 24GB DDR3-1333
- **128-bit vector width**
- **Up to SSE 4.2 ISA**
- 45nm
- Q2 2009

Haswell (HSW)

- Core i7 5930K
- 3.5GHz hex-core
- 15MB L3 cache
- 32GB DDR4-2133
- **256-bit vector width**
- **Up to AVX2 ISA**
- 22nm
- Q4 2013

Processor dynamic frequency scaling was disabled for all experiments
We only analyzed speedup for cross-architecture comparison, not performance

Execution Configuration

- Each loop nest executed 10 times
- Ensured execution times within 1% (0.8% observed)
- Verified correctness of execution for each version by dumping live-out data (PolyBench loop nests) or checksum (TSVC loop nests)
- Average speedup recorded for each loop nest and optimization sequence pair
 - Used for exhaustive search performance and speedup predictor

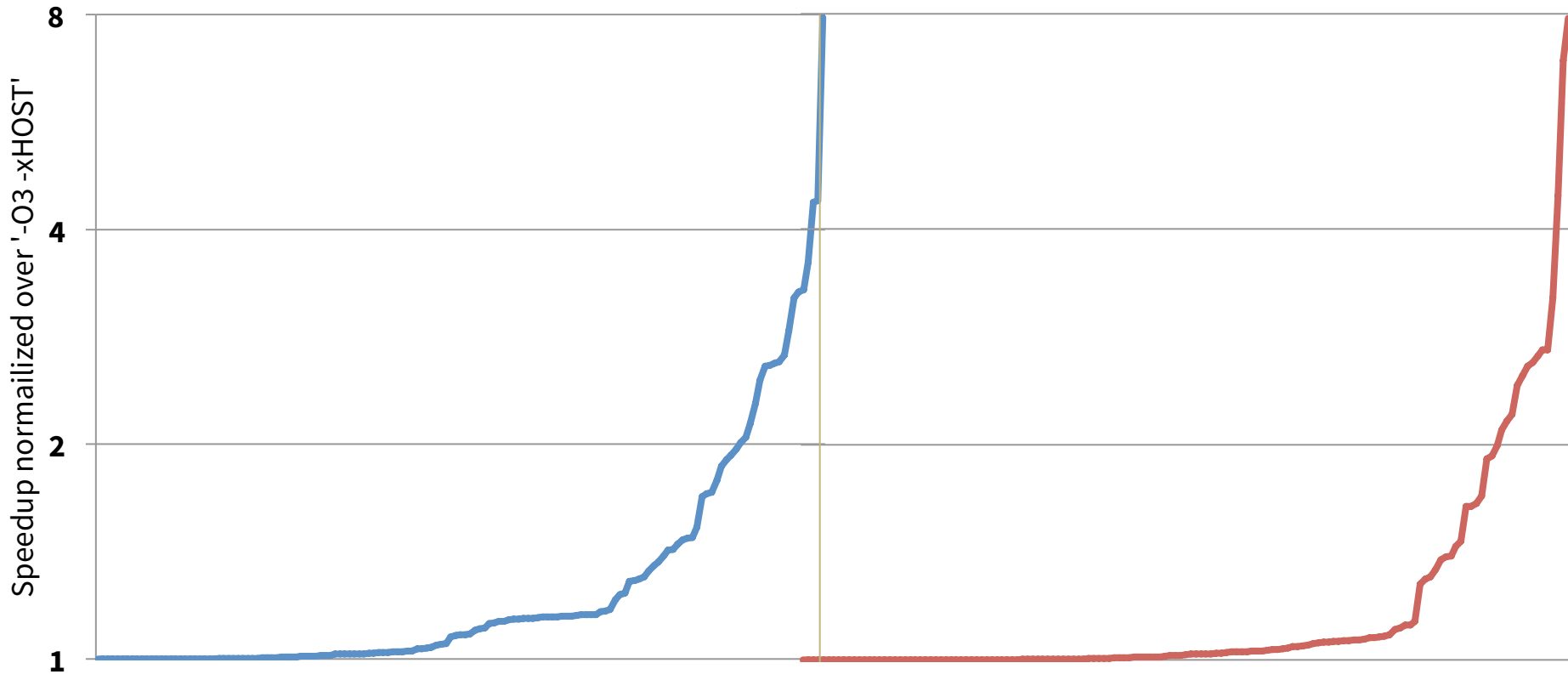
Experiment Results

EXHAUSTIVE SEARCH SPACE SPEEDUP

TSVC Results

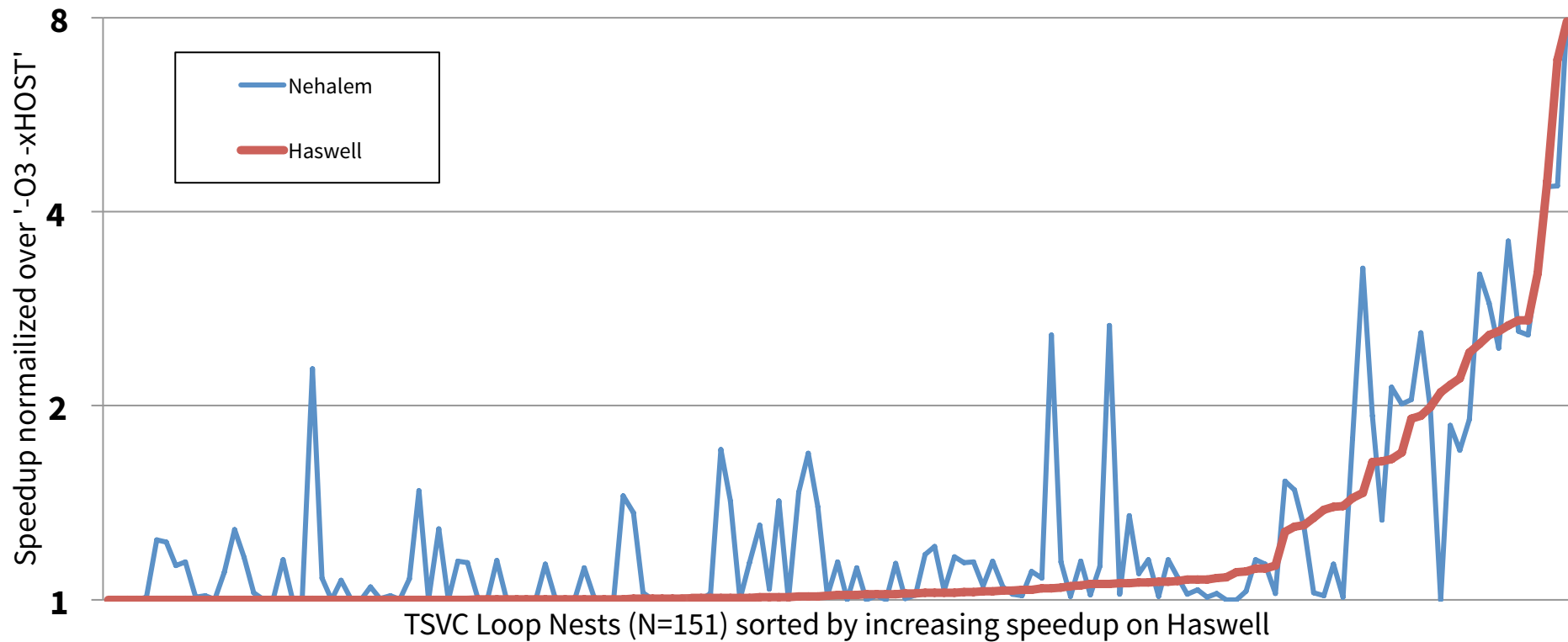
Nehalem

Haswell



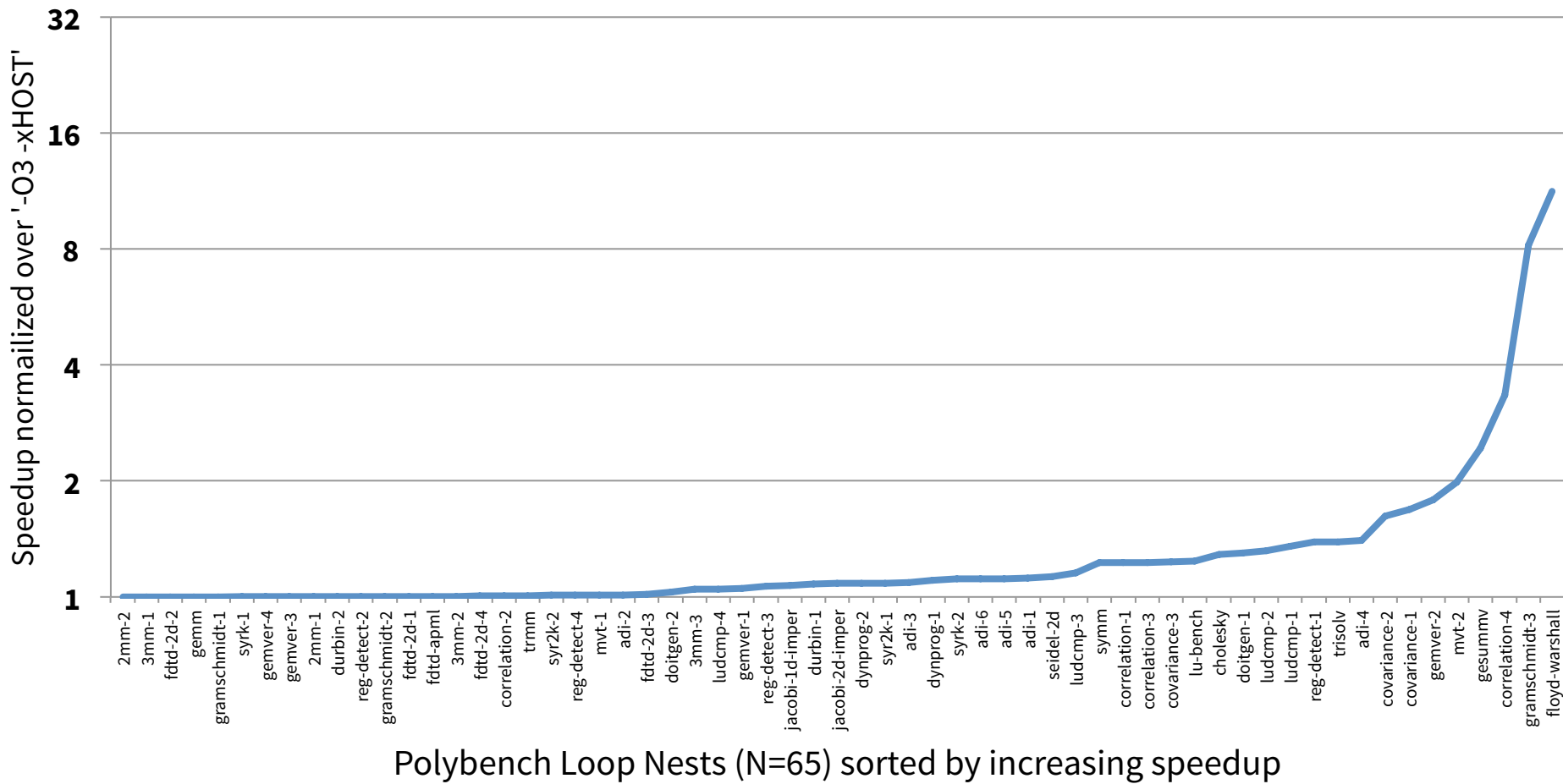
TSVC Loop Nests (N = 151) sorted by increasing speedup

TSVC Cross-Architecture Analysis

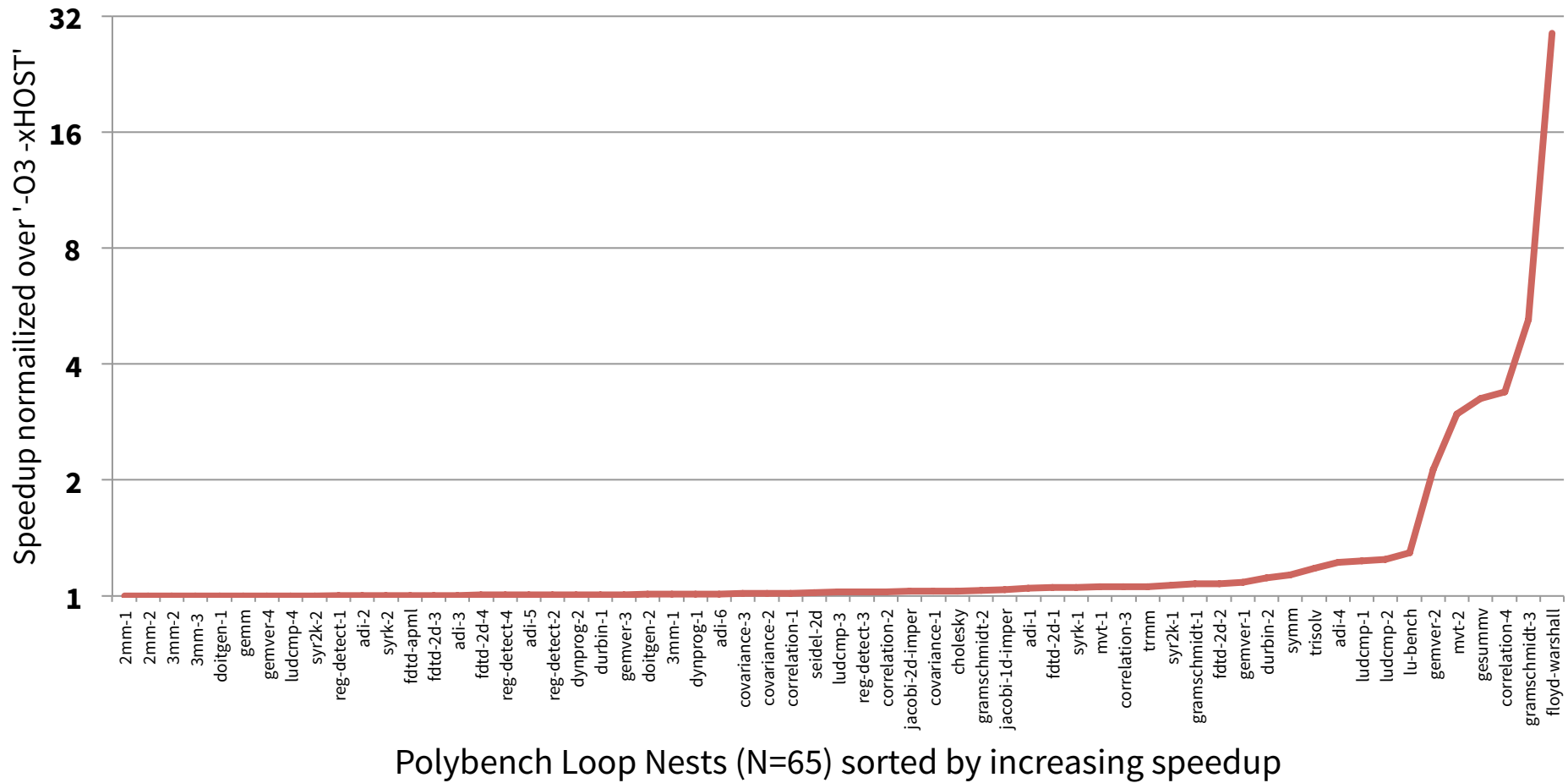


Correlation: $C = 0.8945$

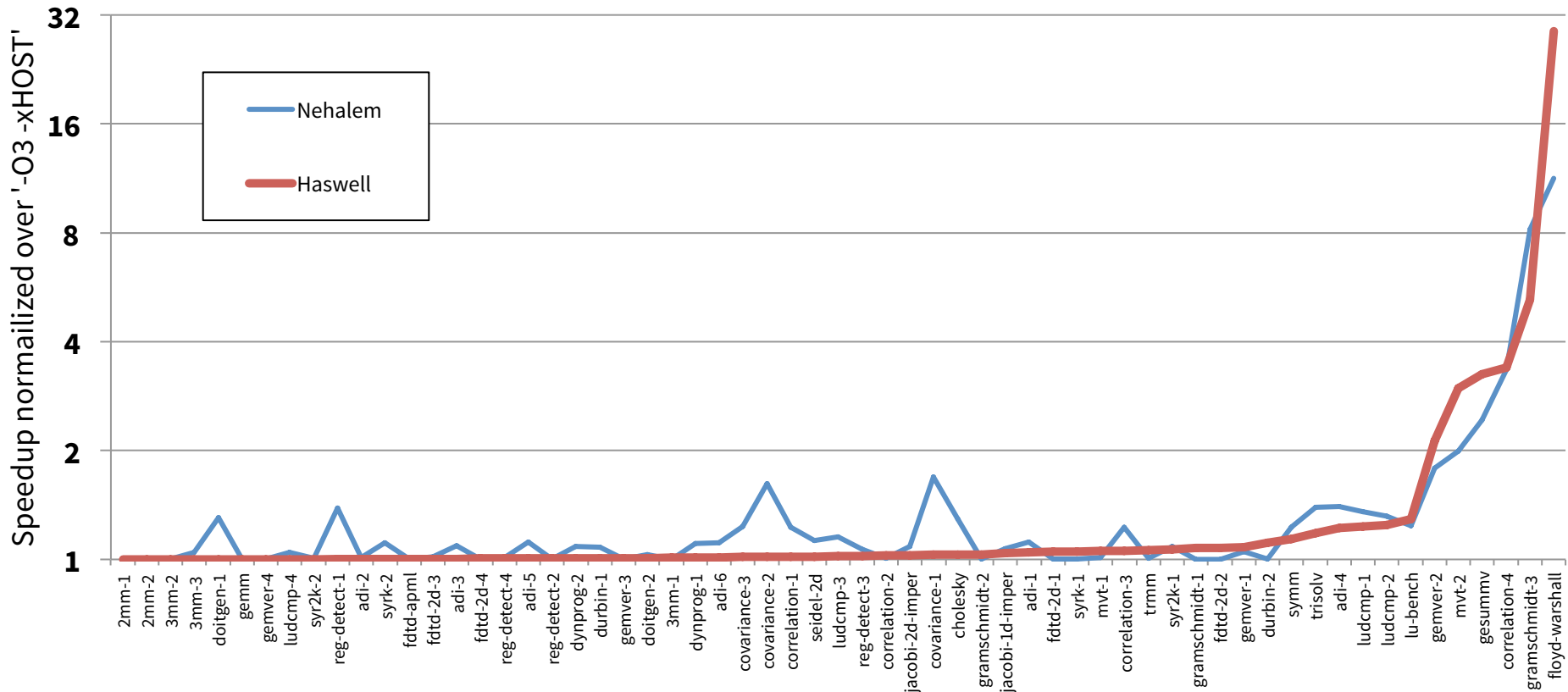
PolyBench Results - Nehalem



PolyBench Results - Haswell



PolyBench Cross-Architecture Analysis



Polybench Loop Nests (N=65) sorted by increasing speedup on Haswell

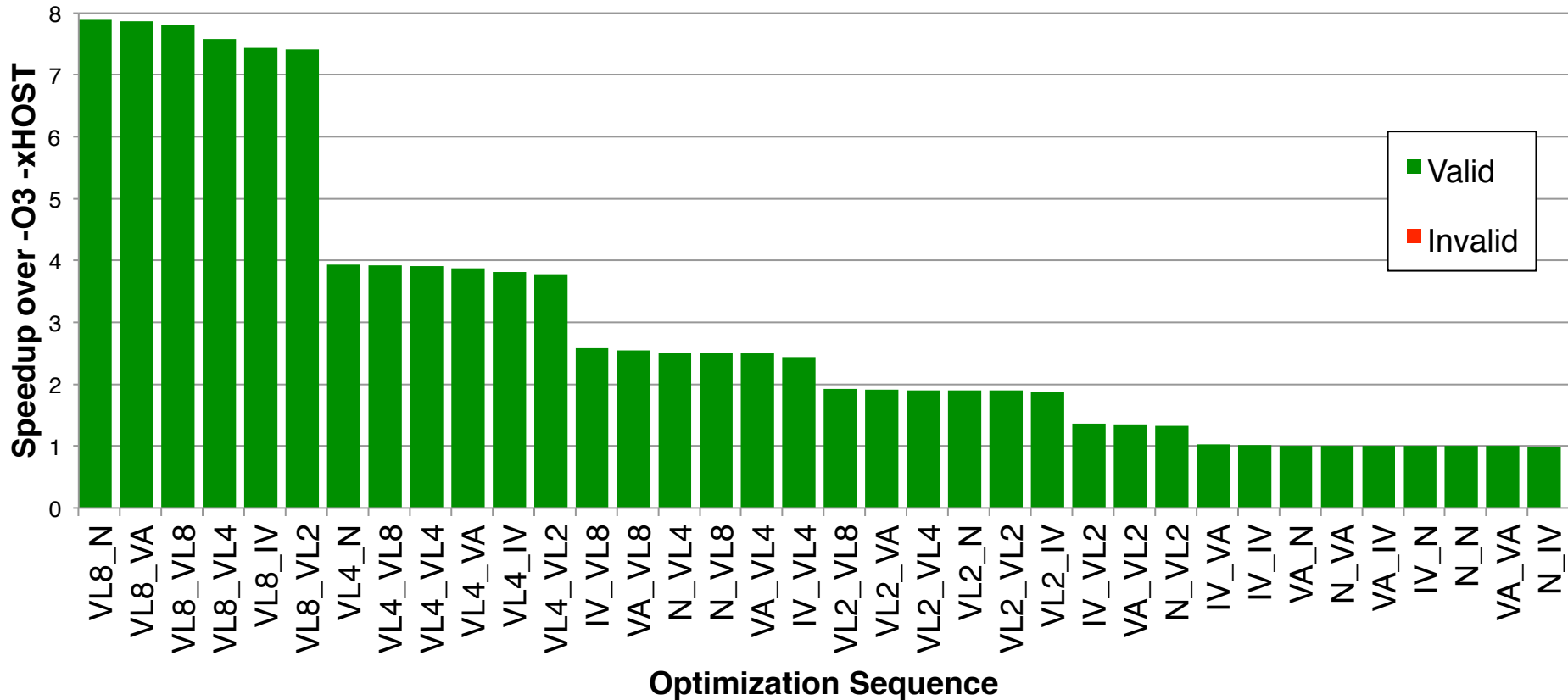
Correlation: C = 0.8894

Experiment Results

EXHAUSTIVE SEARCH SPACE VALID CODE GENERATION

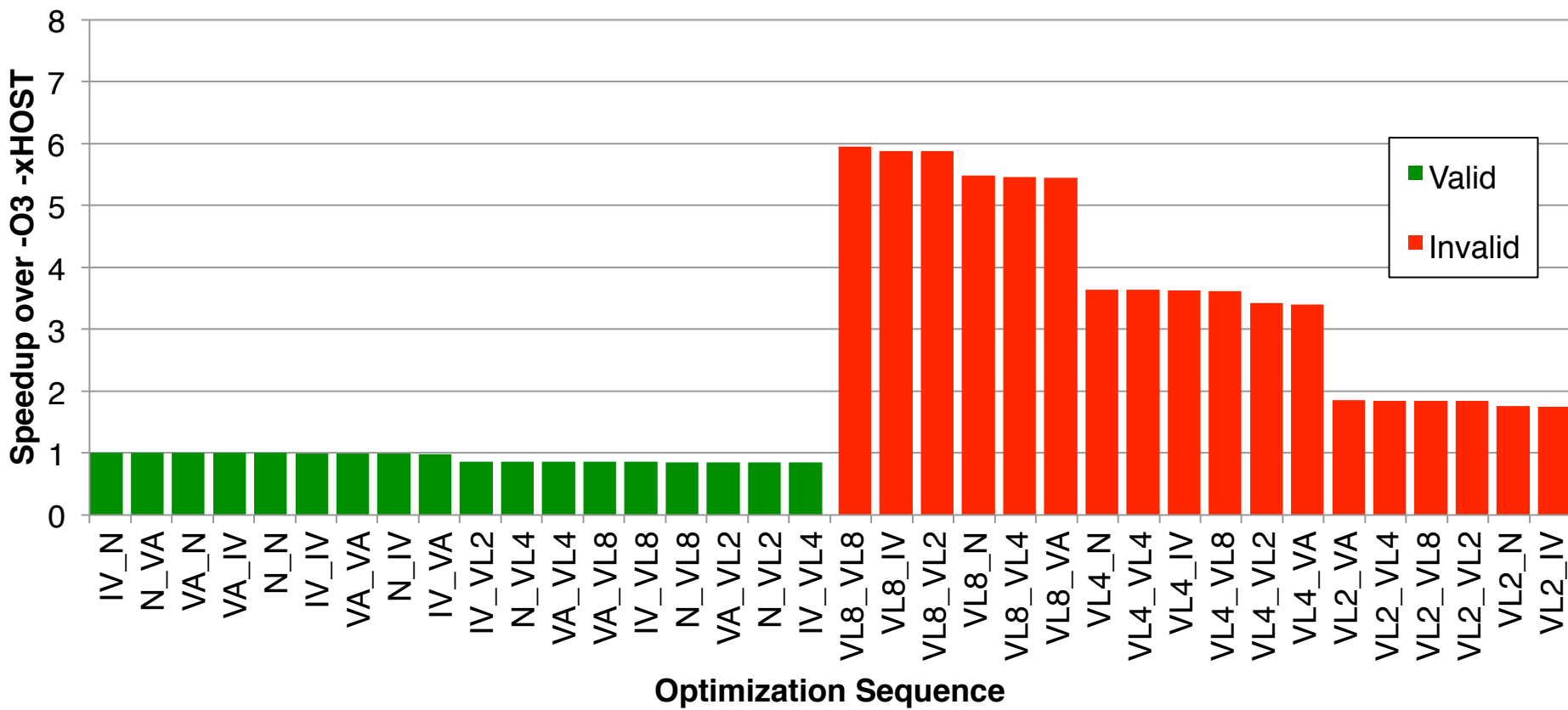
Version Generation Example

T SVC s256 Speedups



Version Generation Example

TSVC s126 Speedups



Version Generation Statistics

Benchmark	Arch	Valid	Invalid	Error
TSVC	Nehalem	1832	151	3
TSVC	Haswell	1826	155	5
PolyBench	Nehalem	5204	3826	0
Polybench	Haswell	5204	3826	0

Version Generation Across Architecture

Architecture	Valid Fastest	Invalid Fastest
Nehalem	140	11
Haswell	143	8

Version Generation Performance for TSVC

Experiment Results

GRAPH-BASED SPEEDUP PREDICTOR

Evaluation Model

- Leave-One-Out Cross Validation
 - For a given loop nest, construct a model based on all other loop nests as training data
 - Compare predictor's speedup to actual speedup
 - 151 models for TSVC, 65 models for PolyBench
- Evaluation Method
 - 1-shot : only consider top prediction
 - 3-shot : consider top three predictions
 - Top prediction is the optimization with best observed speedup

TSVC Speedup Predictor - Nehalem

loop	1-shot	3-shot	Optimal
s126	1.91	1.92	6.28
s221	0.39	1.42	1.42
s2251	2.03	2.44	2.44
s244	0.46	1.36	1.36
s256	0.98	0.99	7.92
s3112	1.99	4.03	4.03
s321	0.50	0.97	2.13
s424	0.99	1.94	2.88
<i>Arith. Mean</i>	<i>0.70</i>	<i>0.98</i>	<i>1.47</i>
<i>Geo. Mean</i>	<i>0.61</i>	<i>0.85</i>	<i>1.32</i>

74% of Optimal on Nehalem 3-shot

TSVC Speedup Predictor - Haswell

loop	1-shot	3-shot	Optimal
s126	1.74	1.84	5.95
s221	0.36	1.35	1.35
s2251	1.49	1.64	1.64
s244	0.32	1.18	1.30
s256	0.99	1.00	7.88
s3112	1.99	4.52	4.52
s321	0.44	1.00	1.65
s424	1.00	1.77	2.57
<i>Arith. Mean</i>	<i>0.62</i>	<i>0.94</i>	<i>1.36</i>
<i>Geo. Mean</i>	<i>0.51</i>	<i>0.80</i>	<i>1.21</i>

77% of Optimal on Haswell 3-shot

PolyBench Speedup Predictor - Nehalem

loop	1-shot	3-shot	Optimal
2mm-1	0.25	0.25	1.00
adi-4	0.99	1.32	1.40
correlation-1	1.00	1.22	1.22
covariance-1	1.37	1.68	1.68
dynprog-1	0.98	1.00	1.10
floyd-warshall	5.22	8.52	11.30
gemm	0.14	0.32	1.00
grammschmidt-3	1.00	7.38	8.17
<i>Arith. Mean</i>	<i>0.97</i>	<i>1.21</i>	<i>1.46</i>
<i>Geo. Mean</i>	<i>0.85</i>	<i>0.98</i>	<i>1.24</i>

84.44% of Optimal on Nehalem 3-shot

PolyBench Speedup Predictor - Haswell

loop	1-shot	3-shot	Optimal
2mm-1	0.67	0.68	1.00
adi-4	1.22	1.22	1.22
correlation-1	0.98	1.00	1.01
covariance-1	1.02	1.02	1.02
dynprog-1	0.99	1.00	1.00
floyd-warshall	25.88	25.88	28.86
gemm	0.12	0.38	1.00
grammschmidt-3	3.37	3.40	5.22
<i>Arith. Mean</i>	<i>1.34</i>	<i>1.47</i>	<i>1.66</i>
<i>Geo. Mean</i>	<i>0.90</i>	<i>1.03</i>	<i>1.20</i>

88.74% of Optimal on Haswell 3-shot

Threats to Validity

- Correctness of generated code
 - PolyBench - Analyzing live out data still may not verify correctness
 - TSVC - used a checksum computation. Invalid results still possible
- Speedup measurement
 - Execution performed on single user mode with timing at a kernel level
 - Speedup a “trend” for PolyBench – not representative of observable speedup for entire kernel
- Machine learning model
 - Optimization bit vector only defines a “level” of vectorization
 - SVM training parameters not explored
 - Generated model seems to find smaller variations between code – similar kernel matrices are generated

Contributions

- VALT directive compiler to simple code generation across different compiler backends
- autovec - exhaustive search code version generator
- Graph-based speedup predictor designed to predict the best vector optimizations to apply to a given loop nest
- Performance analysis of vectorizable micro-benchmarks that can carry across to similar types of kernels

Differences from Related Work

- Stock et al. work only targeted Tensor Contradiction and stencil kernels and didn't use graph-based learning
 - Our approach works on many different types of code and uses graph-based features to construct the model
- Park et al. focused on a different set of optimizations, primarily targeting loop transformations, autoparallelization, and choosing whether or not to vectorize
 - Our approach explores the vectorization search space of loop nests, and allows us to potentially reach a more local maximum speedup given our optimization search space

Future Work

- Extend VALT to support multiple backends (PGI Compiler)
- Change how optimizations are represented
 - Annotate graph-based representation
 - Would eliminate encoding for maximum loop nest size
- Extend work to additional compilers
 - Newer versions of GCC (4.9+), PGI compiler
- Target wider vector size architectures
 - Xeon Phi (Knight's Corner) – 512-bit vector width; limited ISA
 - Knight's Landing and Skylake – AVX-512

Conclusion

- Provided automated and manual techniques for improving performance codes with vectorization optimizations
- Non-experts can use the utilities developed to automatically optimize codes to exploit vector hardware
- With the contributions presented, we
 - achieved up to a 30x speedup through exhaustive search
 - predicted within 88% of search space optimal using the proposed speedup predictor

Thank You!

QUESTIONS?