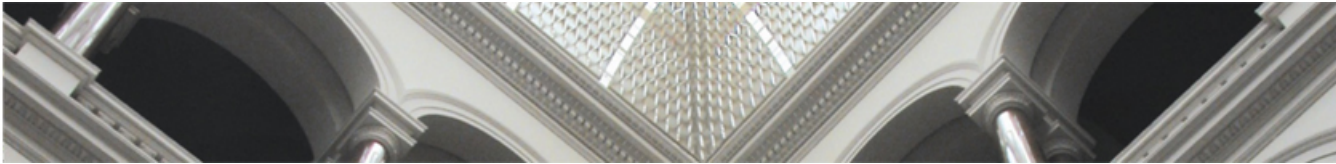




Cost Modelling for Vectorization on ARM

Angela Pohl, Biagio Cosenza and Ben Juurlink

ARM Research Summit 2018



Challenges of Auto-Vectorization in Compilers

1. Is it **possible** to vectorize the code?

- Passes: Loop Level Vectorization (LLV), Superword-Level Parallelism (SLP) ...
- Algorithms: loop transformations, partial vectorization, SLP padding ...
- Challenges: dependences, missing instructions ...

2. Is it **beneficial** to vectorize the code?

- Transformations might add overhead
- Code can run slower afterwards



Short Introduction to Cost Modelling

```
int s121() {
    for (int i = 0; i < LEN; i++){
        j = i + 1;
        a[i] = a[j] + b[i];
    }
    return 0;
}
```

C Source

```
%indvars.iv.next = add nuw nsw i64 %indvars.iv, 1, !dbg !55
%1 = load float, float* %arrayidx, align 4, !dbg !56, !tbaa !22
%2 = load float, float* %arrayidx7, align 4, !dbg !57, !tbaa !22
%add8 = fadd fast float %2, %1, !dbg !58
store float %add8, float* %arrayidx10, align 4, !dbg !60, !tbaa !22
%exitcond = icmp eq i64 %indvars.iv.next, 31999, !dbg !61
```

LLVM's Intermediate Representation

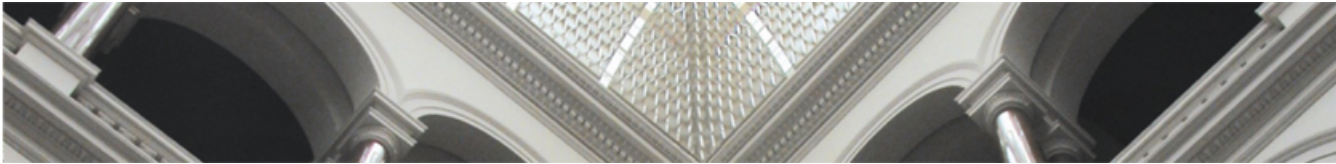
$$S_{est} = \frac{C_{scalar}}{C_{vec}}$$

Estimate Benefit

$$C_{loop} = \frac{1}{VF} \sum C_{instr, VF}$$

```
Cadd, VF=1 = 1
Cload, VF=1 = 1
Cload, VF=1 = 1
Cfadd, VF=1 = 2
Cstore, VF=1 = 1
Cicmp, VF=1 = 1
```

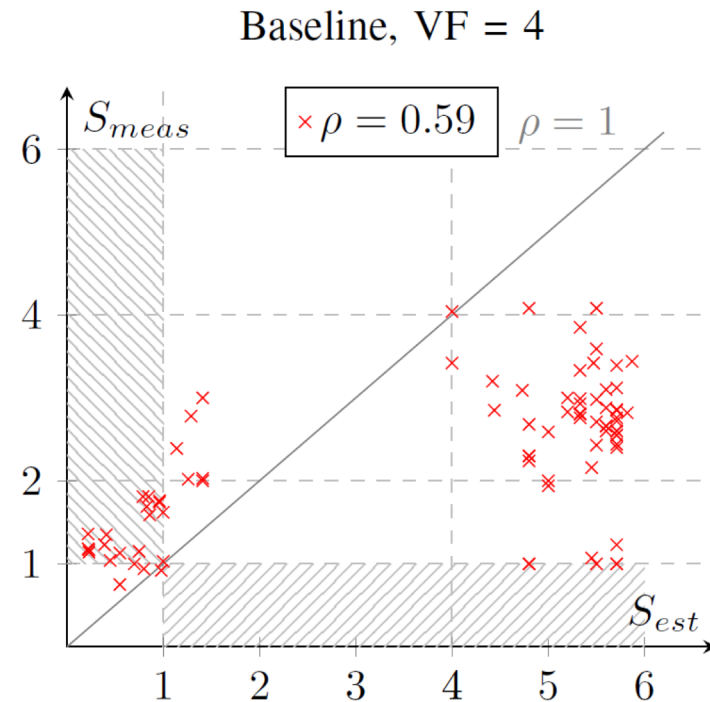
Block Cost Calculation

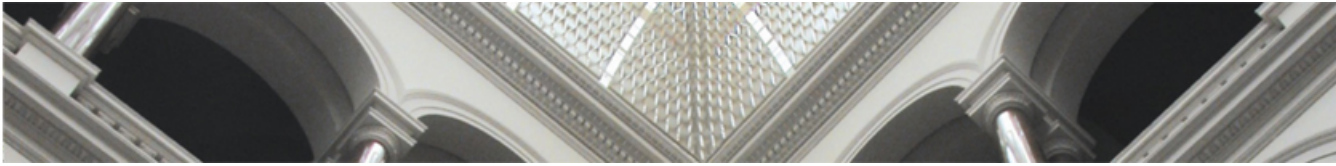


State of the Art Analysis

- LLVM pass of LLVM 6.0 on ARMv8
- TSVC Benchmark: 151 basic loop patterns
- Vectorization with overwritten cost model
- No unrolling, no interleaving

Metric	Baseline
Vectorized loops	61
f_{\ominus}	14
f_{\oplus}	0
t_{norm}	117.2





Linear Modelling Approach

- Formulate basic block as linear equation for each test pattern

$$C_{vec} = \sum n_i C_i$$

- Determine basic block target cost based on measurement

$$C_{target} = \frac{C_{scalar}}{S_{meas}} = \sum n_i C_i'$$

- Apply mathematical fitting to set of linear equations

Linear Modelling: Example

```
for (int i = 0; i < LEN; i++){
    a[i] = b[k] + 1.0;
}
```

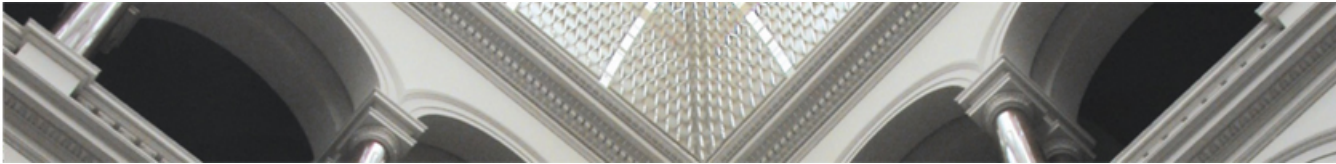
$$C_{scalar} = 8 \quad S_{meas} = 2.76 \Rightarrow C_{target} = 2.89$$

```
for (int i = 0; i < LEN; i++){
    prod *= a[i];
}
```

$$C_{scalar} = 6 \quad S_{meas} = 2.30 \Rightarrow C_{target} = 2.61$$

$$2.89 = 1 \times C_{load} + 1 \times C_{fadd} + 0 \times C_{fmult} + \dots$$

$$2.61 = 1 \times C_{load} + 0 \times C_{fadd} + 1 \times C_{fmult} + \dots$$



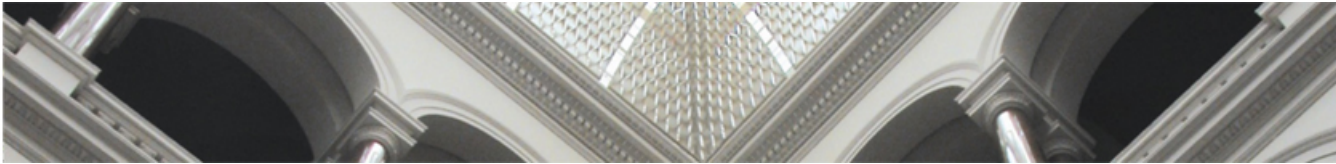
Modelling Speedup Instead of Cost

Problem:

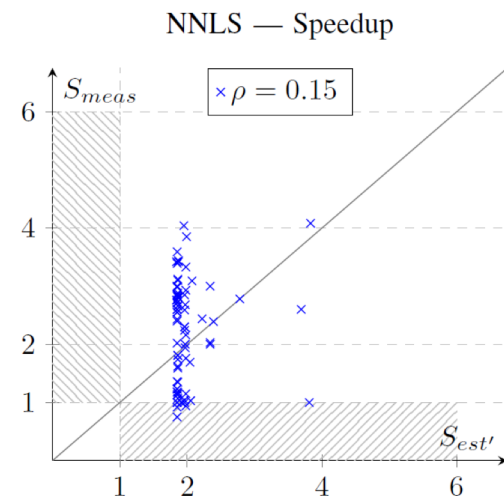
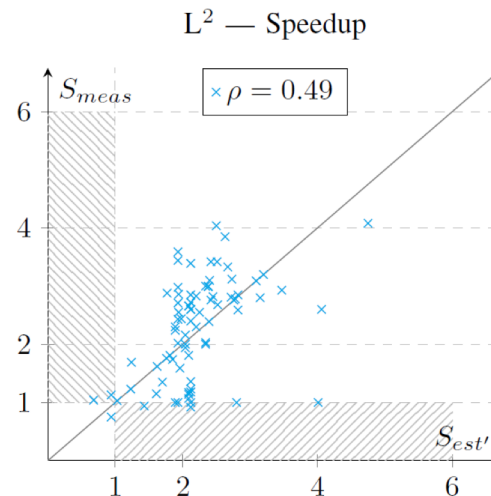
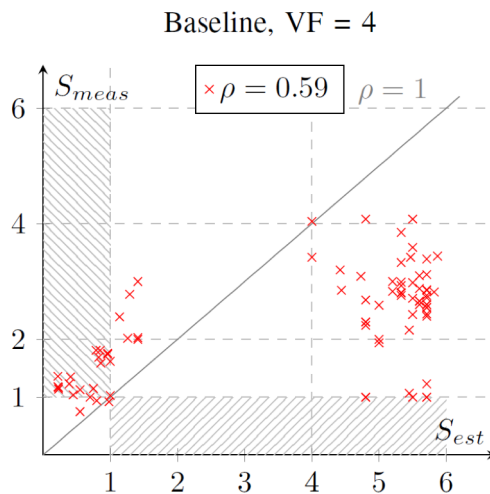
- Target costs vary significantly, i.e. data has to be fitted across large interval
- But: fitting benefits from smaller target intervals

Idea:

- Model speedup instead of cost
- Limits interval to (0, vectorization factor)
- Linear equations become $S_{meas} = \sum n_i w_i$
with n_i = numer of instructions of same type, w_i = corresponding weight



Results: Fitted for Speedup



- L2: Least Squares (minimizes Euclidian L2 Norm)
- NNLS: Non-Negative Least Squares (all coefficients > 0)



Adding Block Composition as Feature

Problem:

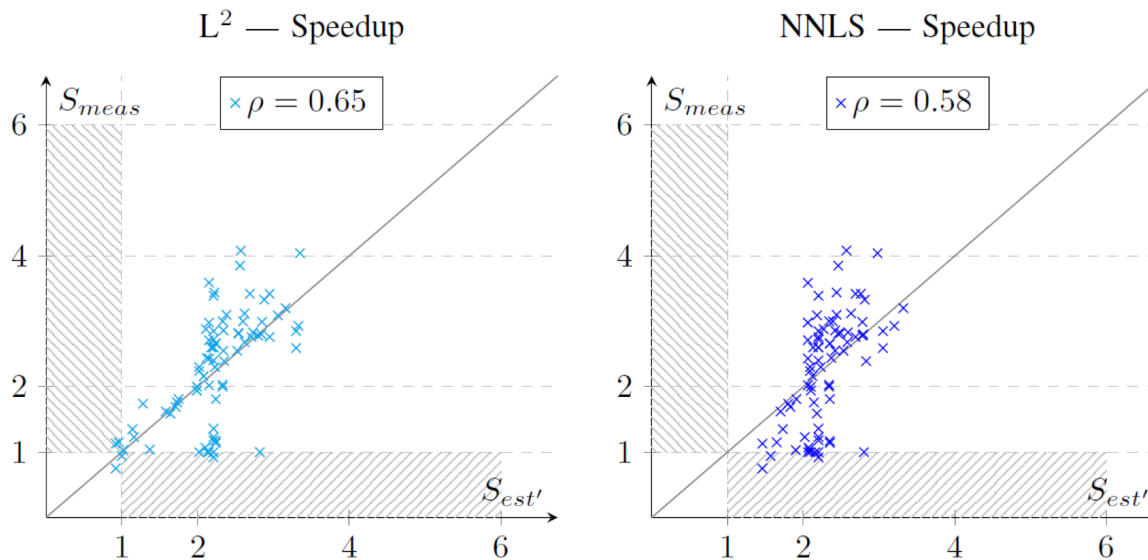
- Current cost model only cares about individual instructions
- Arithmetic intensity can have major impact on speedup, e.g. if code is memory bound

Idea:

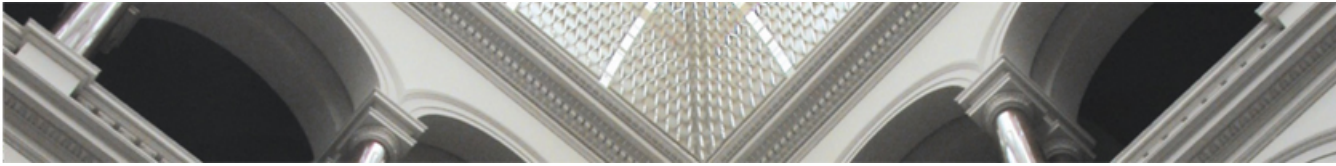
- Replace count of instruction type with overall percentage, e.g. 20% load, 10% cmp, ...

$$S_{meas} = \sum \frac{n_i}{n_{total}} w_i$$

Results: Fitted with Rated Instruction Count

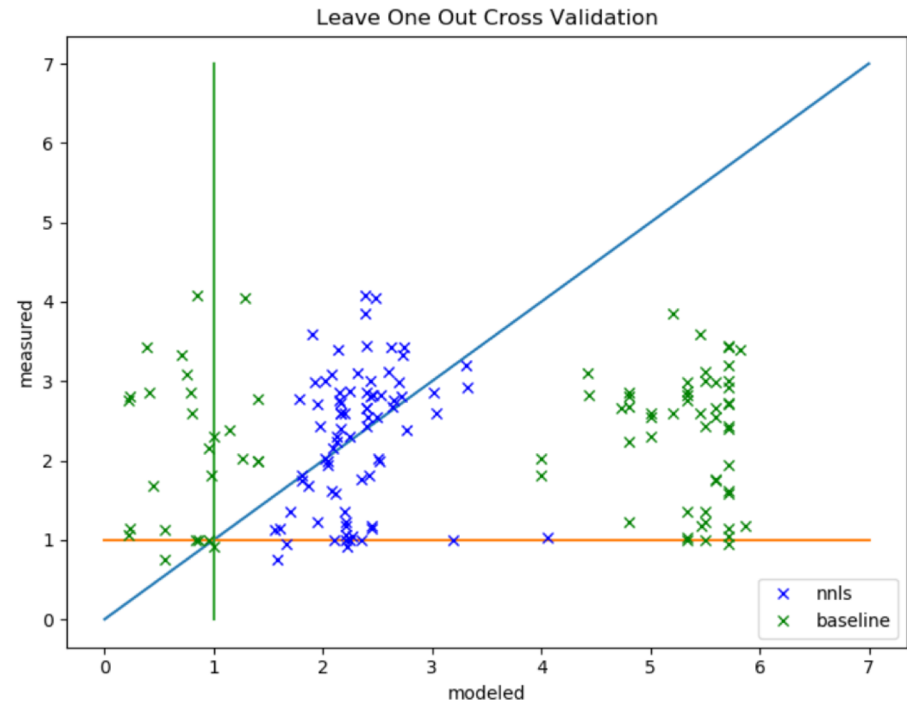


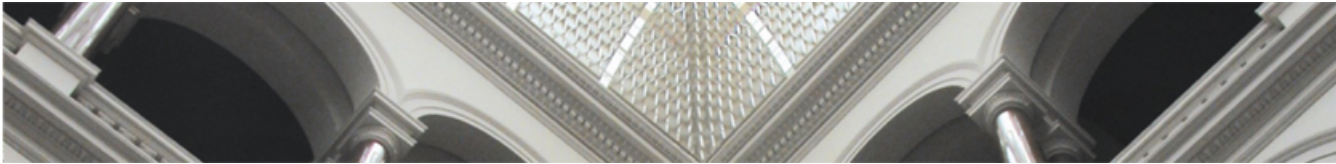
Metric	Baseline	L^2	NNLS
Vectorized loops	61	76	80
f_{\ominus}	14	2	0
f_{\oplus}	0	1	3
t_{norm}	117.2	113.2	113.4



Leave One Out Cross Validation: NNLS

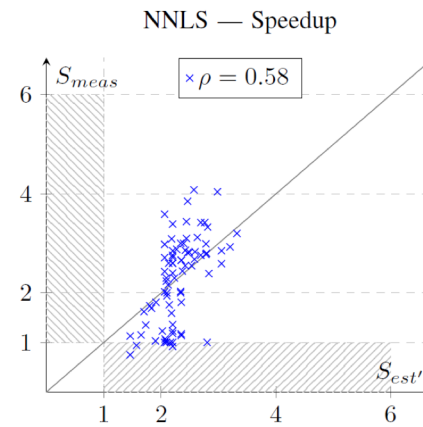
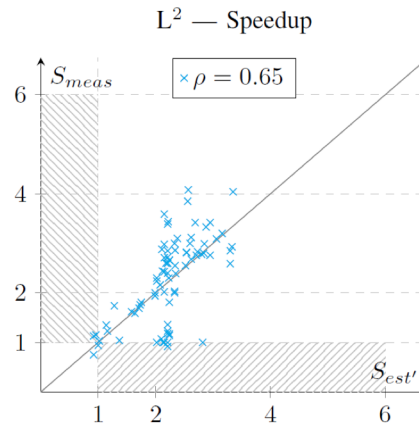
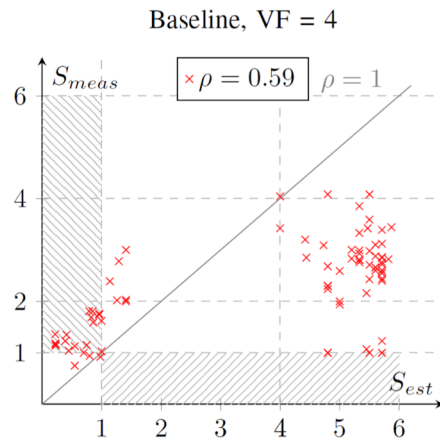
- Fitting is done with whole training data set except for one kernel
- Kernel is then predicted with that model



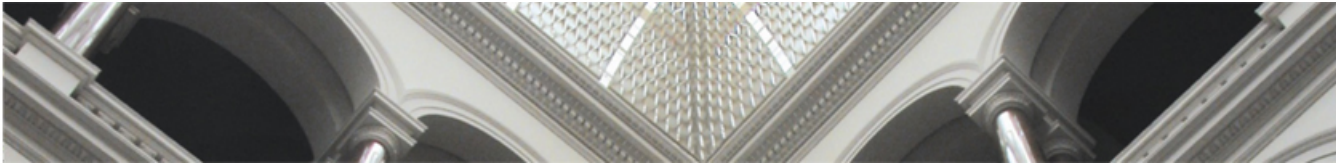


Conclusion

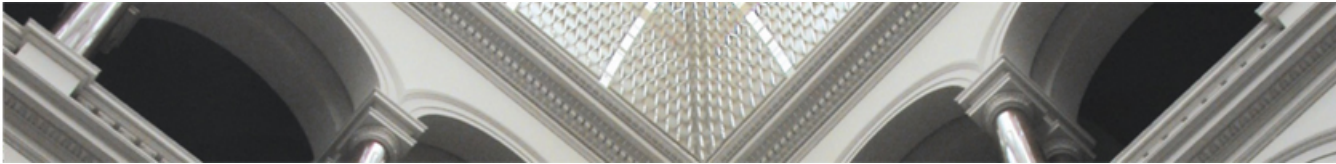
- Compilers need more accurate cost models to avoid mispredictions
- Aligned cost models enable comparison of different transformation options
- With our refined cost model we
 - Increase the correlation between estimated and measured speedup
 - Decrease the number of false predictions
 - Lower execution times
- Next steps: add more code features and tests to cover all instruction types



THANK YOU



BACKUP



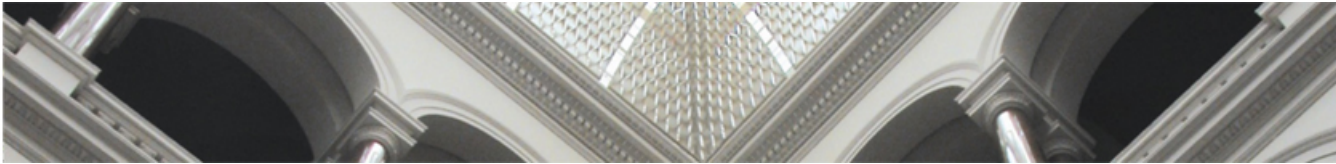
Why a More Accurate Cost Model?

1. Improve classification results to decide whether code should be vectorized or not.
2. Enable comparison of different transformations **with each other**, not just to baseline.

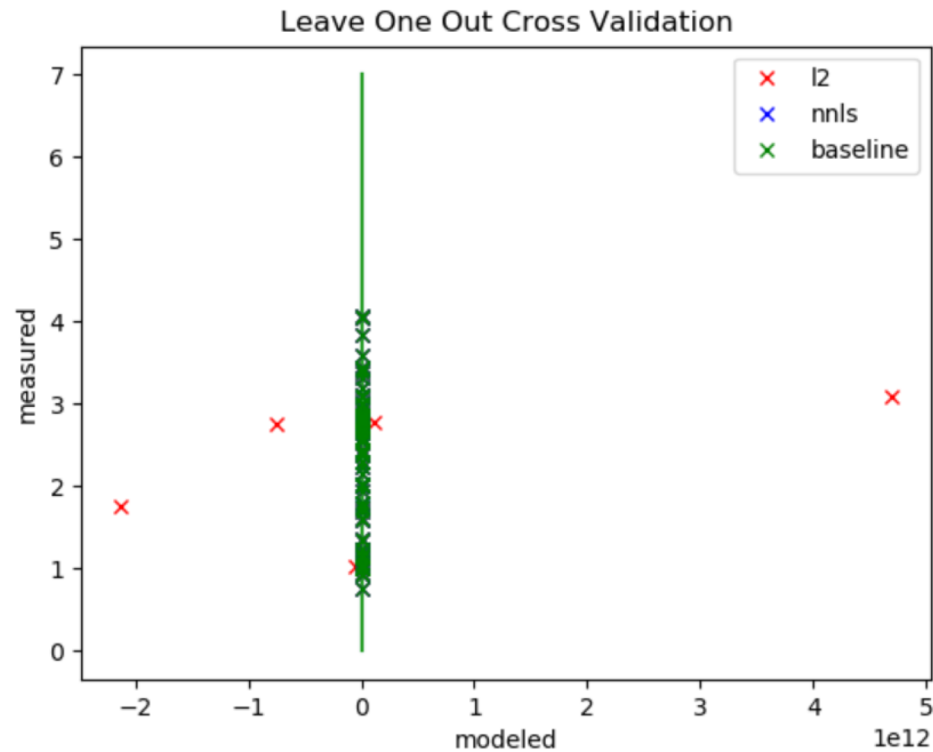
```
for (int i = 0; i < LEN/2; i++){  
    k = j + 1;  
    a[i] = b[k] - d[i];  
    j = k + 1  
    b[k] = a[i] + c[k];  
}
```

	LLV	SLP
Predicted Speedup	0.96	1.00
Measured Speedup	1.13	1.20

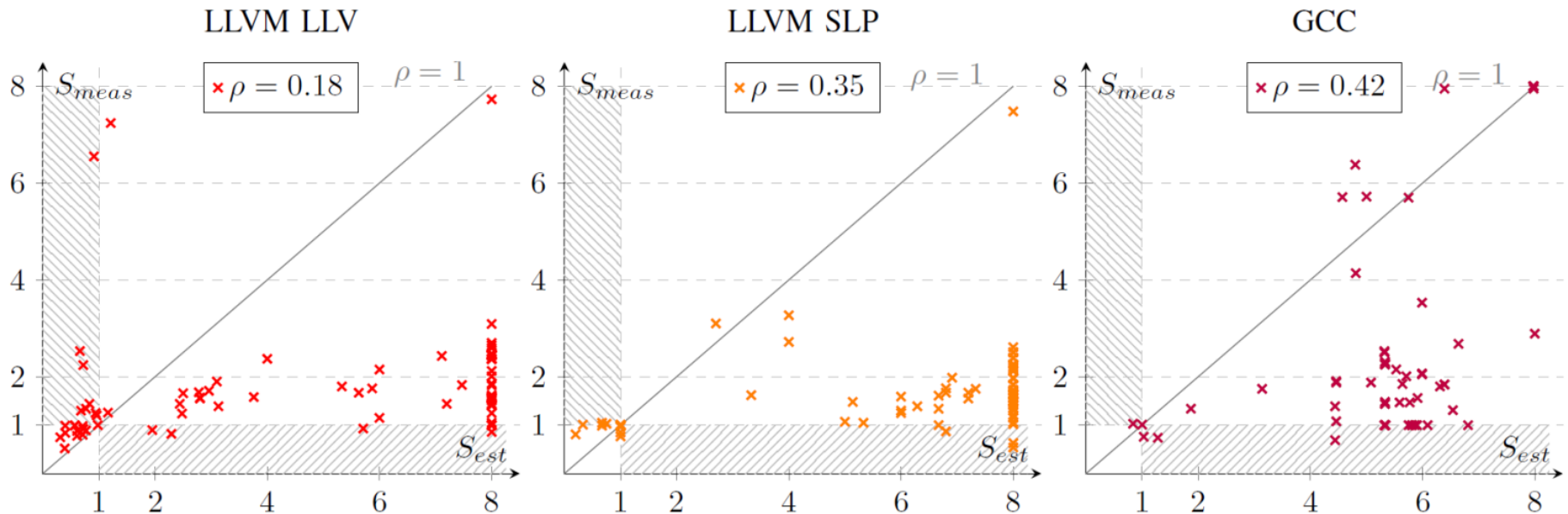
Example code run on Intel i5



L2- LOOCV Validation Results

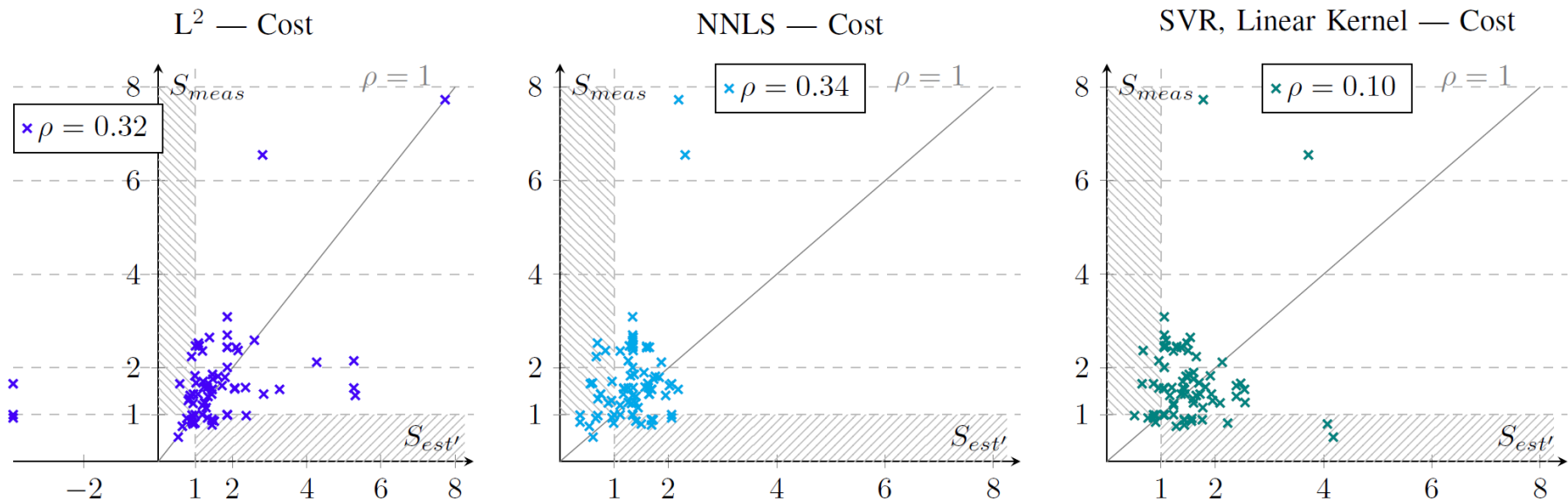


State of the Art Analysis – x86



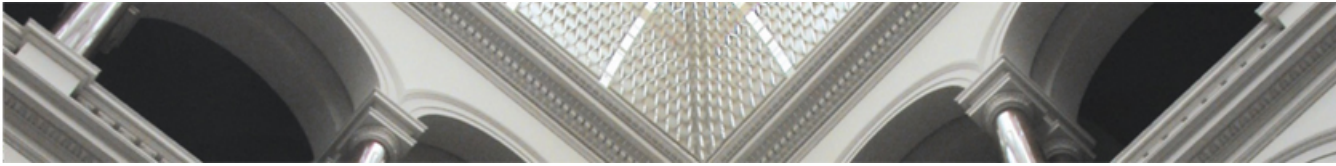
- TSVC benchmark: 151 basic loop patterns
- SLP vectorization applied after loop unrolling
- Results shown for Intel Xeon E5 with AVX2

Results: Fitted for Cost – x86

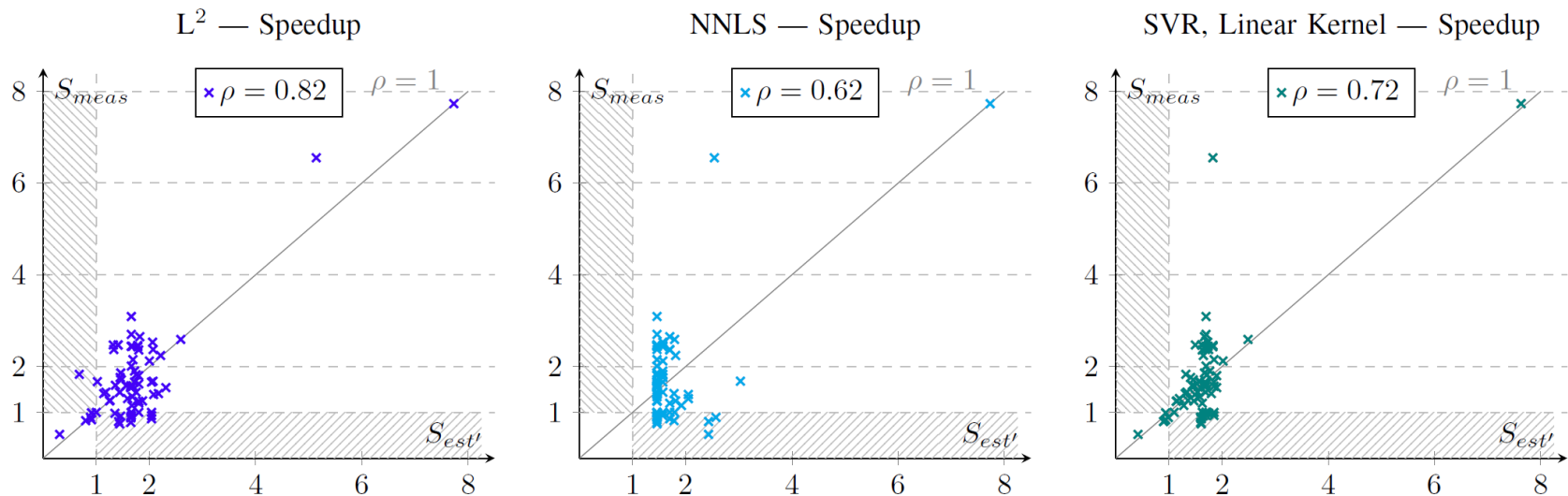


- Correlation between estimated and measured speedup in LLVM's LLV pass after fitting
- L2: Least Squares (minimizes Euclidian L2 Norm)
- NNLS: Non-Negative Least Squares (all coefficients > 0)
- SVR: Support Vector Regression

Cost Modelling for Vectorization on ARM, A. Pohl et al.



Results: Fitted for Speedup – x86



- All three approaches improve correlation further
- False negatives reduced (L2) or eliminated (NNLS, SVR)
- But: small increase in false positives