# SLOTHY: Using Constraint-Solving for Superoptimization of Cryptographic Assembly

Tutorial (with assignment solutions)

Amin Abdulrahman, Max Planck Institute for Security and Privacy, Germany
Matthias J. Kannwischer, Chelpis Quantum Corp., Taiwan
September 14, 2025

# Tutorial Structure

1. Introduction
2. SLOTHY basics
3. Assignment 1 & 2: Basic use of SLOTHY

4. Heuristics and register spilling
5. Assignment 3: Optimizing a large piece of code (Keccak)

6. SLOTHY's Architecture & Microarchitecture model
7. Advanced SLOTHY features
8. Assignment 4 & 5: Extending SLOTHY

# Slides & Assignments



Tutorial Assignments
github.com/dop-amin/ches2025-slothy-tutorial



Tutorial Slides
kannwischer.eu/talks/20250914_slothy.pdf

# Motivation

*"Hey, we need a fast ML-KEM implementation for our new smartphone CPU. Can you do that?"*

*"Hey, we need a fast ML-KEM implementation for our new smartphone CPU. Can you do that?"*

*"Should be easy, right? I'll just use the C reference and take it from there"*

*"Hey, we need a fast ML-KEM implementation for our new smartphone CPU. Can you do that?"*

*"Should be easy, right? I'll just use the C reference and take it from there"*

Narrator: It was not easy.

# Why Cryptographic Engineering is (Not) Easy

- Simplicity

```
1  void ntt(int16_t r[256]) {
2    unsigned int len, start, j, k;
3    int16_t t, zeta;
4
5    k = 1;
6    for(len = 128; len >= 2; len >>= 1) {
7      for(start = 0; start < 256; start = j + len) {
8        zeta = zetas[k++];
9        for(j = start; j < start + len; j++) {
10         t = fqmul(zeta, r[j + len]);
11         r[j + len] = r[j] - t;
12         r[j] = r[j] + t;
13       }
14     }
15   }
16 }
```

- Simplicity
- Security

```
1  for (i = 0; i < KYBER_N / 8; i++){
2    for (j = 0; j < 8; j++){
3      mask = -(int16_t)((msg[i] >> j) & 1);
4      r->coeffs[8 * i + j] = mask & ((KYBER_Q + 1) / 2);
5    }
6  }
```

- Simplicity
- Security
- Performance

```
1  .macro mulmodq dst, src, const, idx0, idx1
2    vqrdmulhq   t2,  \src, \const, \idx1
3    vmulq       \dst, \src, \const, \idx0
4    vmlaq       \dst,  t2, consts, 0
5  .endm
```

- Simplicity
- Security
- Performance
- More Performance

```
1   mul v8.8H, v8.8H, v1.H[4]
2   add v3.8H, v3.8H, v9.8H
3   sub v9.8H, v12.8H, v10.8H
4   add v12.8H, v12.8H, v10.8H
5   mls v8.8H, v16.8H, v7.H[0]
6   str q3, [x0], #(16)
7   ldr q10, [x0, #0]
8   sub v3.8H, v18.8H, v8.8H
9   add v16.8H, v18.8H, v8.8H
10  str q23, [x0, #48]
11  ldr q5, [x0, #64]
12  str q12, [x0, #112]
```

- Simplicity
- Security
- Performance
- More Performance
- Effort

```
1    add  v18.8H, v8.8H, v22.8H
2    mul  v27.8H, v31.8H, v0.H[0]
3    add  v22.8H, v16.8H, v11.8H
4    sqrdmulh v14.8H, v31.8H, v0.H[1]
5    str  q4, [x0, #304]
6    str  q18, [x0, #240]
7    sub  v4.8H, v10.8H, v5.8H
8    add  v18.8H, v10.8H, v5.8H
9    sqrdmulh v24.8H, v22.8H, v0.H[3]
10   mul  v10.8H, v22.8H, v0.H[2]
11   str  q4, [x0, #432]
12   str  q18, [x0, #368]
```

- Simplicity
- Security
- Performance
- More Performance
- Effort
- Auditability & Maintainability

```
1   mul v8.8H, v8.8H, v1.H[4]
2   add v3.8H, v3.8H, v9.8H
3   sub v9.8H, v12.8H, v10.8H
4   add v12.8H, v12.8H, v10.8H
5   mls v8.8H, v16.8H, v7.H[0]
6   str q3, [x0, #(16)
7   ldr q10, [x0, #0]
8   sub v3.8H, v18.8H, v8.8H
9   add v16.8H, v18.8H, v8.8H
10  str q23, [x0, #48]
11  ldr q5, [x0, #64]
12  str q12, [x0, #112]
```

```
1   add v18.8H, v8.8H, v22.8H
2   mul v27.8H, v31.8H, v0.H[0]
3   add v22.8H, v16.8H, v11.8H
4   sqrdmulh v14.8H, v31.8H, v0.H[1]
5   str q4, [x0, #304]
6   str q18, [x0, #240]
7   sub v4.8H, v10.8H, v5.8H
8   add v18.8H, v10.8H, v5.8H
9   sqrdmulh v24.8H, v22.8H, v0.H[3]
10  mul v10.8H, v22.8H, v0.H[2]
11  str q4, [x0, #432]
12  str q18, [x0, #368]
```

# SLOTHY High-Level Workflow

```
1   .macro mulmodq dst, src, const, idx0, idx1
2     sqrdmulh tmp2.8h, \src.8h, \const.h[\idx1]
3     mul \dst.8h, \src.8h, \const.h[\idx0]
4     mla \dst.8h, tmp2.8h, consts.h[0]
5   .endm
6   .macro ct_butterfly a, b, root, idx0, idx1
7     mulmodq tmp, \b, \root, \idx0, \idx1
8     sub \b.8h, \a.8h, tmp.8h
9     add \a.8h, \a.8h, tmp.8h
10  .endm
11
12  ct_butterfly data0, data8, root0, 0, 1
13  ct_butterfly data1, data9, root0, 0, 1
14  ct_butterfly data2, data10, root0, 0, 1
15  ct_butterfly data3, data11, root0, 0, 1
16  ct_butterfly data4, data12, root0, 0, 1
17  ct_butterfly data5, data13, root0, 0, 1
18  ct_butterfly data6, data14, root0, 0, 1
19  ct_butterfly data7, data15, root0, 0, 1
```
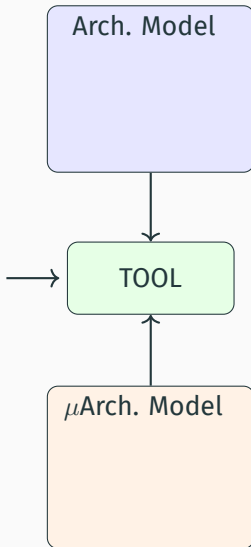
```
1   .macro mulmodq dst, src, const, idx0, idx1
2     sqrdmulh tmp2.8h, \src.8h, \const.h[\idx1]
3     mul \dst.8h, \src.8h, \const.h[\idx0]
4     mla \dst.8h, tmp2.8h, consts.h[0]
5   .endm
6   .macro ct_butterfly a, b, root, idx0, idx1
7     mulmodq tmp, \b, \root, \idx0, \idx1
8     sub \b.8h, \a.8h, tmp.8h
9     add \a.8h, \a.8h, tmp.8h
10  .endm
11
12  ct_butterfly data0, data8, root0, 0, 1
13  ct_butterfly data1, data9, root0, 0, 1
14  ct_butterfly data2, data10, root0, 0, 1
15  ct_butterfly data3, data11, root0, 0, 1
16  ct_butterfly data4, data12, root0, 0, 1
17  ct_butterfly data5, data13, root0, 0, 1
18  ct_butterfly data6, data14, root0, 0, 1
19  ct_butterfly data7, data15, root0, 0, 1
```
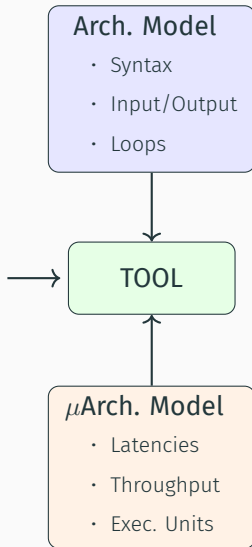
$\longrightarrow$ TOOL

# SLOTHY High-Level Workflow

```
1   .macro mulmodq dst, src, const, idx0, idx1
2     sqrdmulh tmp2.8h, \src.8h, \const.h[\idx1]
3     mul \dst.8h, \src.8h, \const.h[\idx0]
4     mla \dst.8h, tmp2.8h, consts.h[0]
5   .endm
6   .macro ct_butterfly a, b, root, idx0, idx1
7     mulmodq tmp, \b, \root, \idx0, \idx1
8     sub \b.8h, \a.8h, tmp.8h
9     add \a.8h, \a.8h, tmp.8h
10  .endm
11
12  ct_butterfly data0, data8, root0, 0, 1
13  ct_butterfly data1, data9, root0, 0, 1
14  ct_butterfly data2, data10, root0, 0, 1
15  ct_butterfly data3, data11, root0, 0, 1
16  ct_butterfly data4, data12, root0, 0, 1
17  ct_butterfly data5, data13, root0, 0, 1
18  ct_butterfly data6, data14, root0, 0, 1
19  ct_butterfly data7, data15, root0, 0, 1
```



Arch. Model

TOOL

$\mu$Arch. Model

```
1   .macro mulmodq dst, src, const, idx0, idx1
2     sqrdmulh tmp2.8h, \src.8h, \const.h[\idx1]
3     mul \dst.8h, \src.8h, \const.h[\idx0]
4     mla \dst.8h, tmp2.8h, consts.h[0]
5   .endm
6   .macro ct_butterfly a, b, root, idx0, idx1
7     mulmodq tmp, \b, \root, \idx0, \idx1
8     sub \b.8h, \a.8h, tmp.8h
9     add \a.8h, \a.8h, tmp.8h
10  .endm
11
12  ct_butterfly data0, data8, root0, 0, 1
13  ct_butterfly data1, data9, root0, 0, 1
14  ct_butterfly data2, data10, root0, 0, 1
15  ct_butterfly data3, data11, root0, 0, 1
16  ct_butterfly data4, data12, root0, 0, 1
17  ct_butterfly data5, data13, root0, 0, 1
18  ct_butterfly data6, data14, root0, 0, 1
19  ct_butterfly data7, data15, root0, 0, 1
```
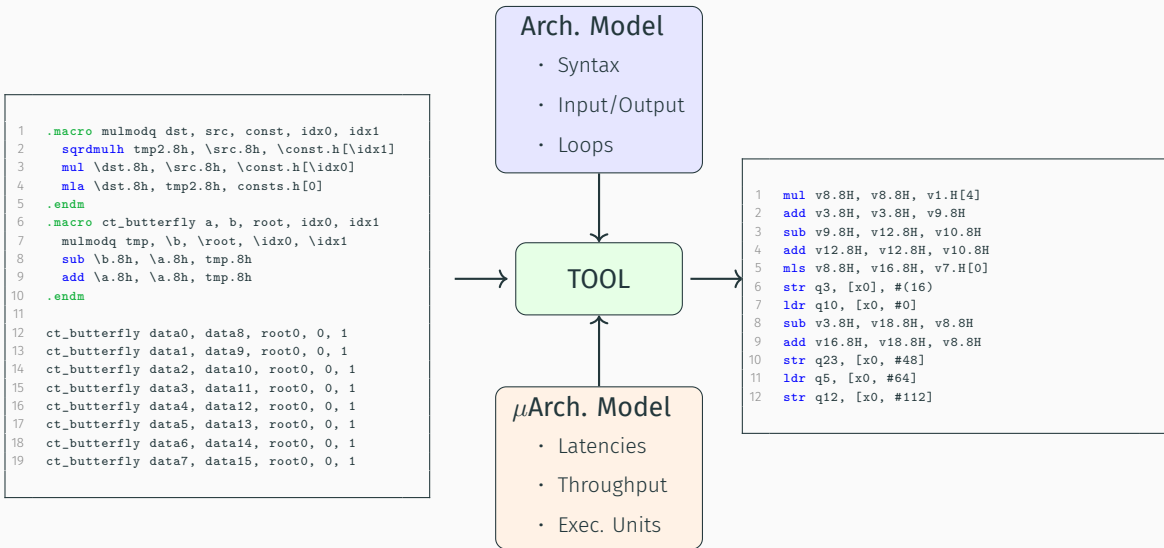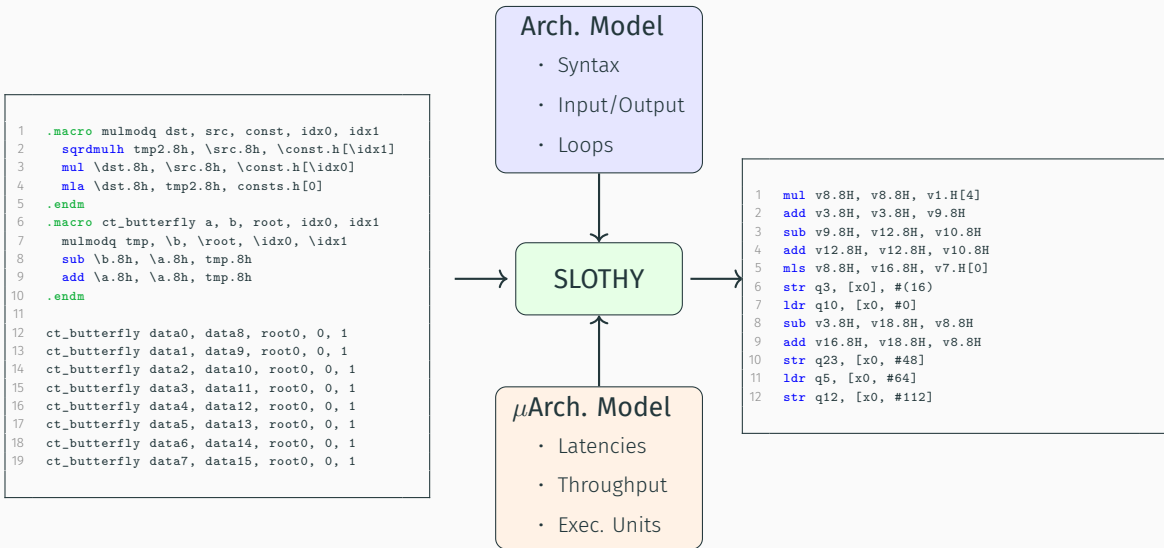
Arch. Model
· Syntax
· Input/Output
· Loops

TOOL

$\mu$Arch. Model
· Latencies
· Throughput
· Exec. Units

5

# SLOTHY High-Level Workflow

# SLOTHY High-Level Workflow



5

SLOTHY: Super (Lazy) Optimization of Tricky Handwritten assemblY

# What is SLOTHY?

SLOTHY: Super (Lazy) Optimization of Tricky Handwritten assemblY

A **fixed-instruction superoptimizer** that:

- Takes your assembly code as input
- Preserves your instruction choices (security!)
- Simultaneously solves
  - instruction scheduling
  - register allocation
  - software pipelining

# What is SLOTHY?

SLOTHY: Super (Lazy) Optimization of Tricky Handwritten assemblY

A **fixed-instruction superoptimizer** that:

- Takes your assembly code as input
- Preserves your instruction choices (security!)
- Simultaneously solves
  - instruction scheduling
  - register allocation
  - software pipelining

### Finding Optimal Solutions

Formulates optimization as a **constraint satisfaction problem** and uses Google OR-Tools' CP-SAT to find optimal solutions.

## Supported Platforms

Partial models for a variety of architectures and microarchitectures:

### AArch64 (+ Neon)

- Arm Cortex-A55
- Arm Cortex-A72
- Arm Neoverse N1
- Apple M1

### Armv8-M (+ MVE)

- Arm Cortex-M55
- Arm Cortex-M85

### Armv7E-M

- Arm Cortex-M7
- Arm Cortex-M4 (WIP)

### RISC-V (WIP)

- XuanTie C908 (WIP)
- SpacemiT X60 (WIP)

| Workload | $\mu$Arch | Before (cycles) | After (cycles) | Speed-up |
| --- | --- | --- | --- | --- |
| ML-DSA NTT | Cortex-A55 | 2436 | 1728 | 1.41× |
| ML-DSA NTT | Cortex-A72 | 2241 | 1766 | 1.27× |
| ML-DSA NTT | Cortex-M7 | 8139 | 4141 | 1.97× |
| X25519 | Cortex-A55 | 143849 | 139752 | 1.03× |
| Keccak-f1600 | Cortex-M7 | 6691 | 5149 | 1.30× |

Also successfully applied to complex FFT, more subroutines from ML-DSA & ML-KEM

A. Abdulrahman, H. Becker, M. J. Kannwischer, and F. Klein. *Fast and Clean: Auditable high-performance assembly via constraint solving.* CHES '24. 2023

A. Abdulrahman, M. J. Kannwischer, and T.-H. Lim. "Enabling Microarchitectural Agility: Taking ML-KEM & ML-DSA from Cortex-M4 to M7 with SLOTHY". In: ASIA CCS '25. 2025

# SLOTHY is Used in Practice

SLOTHY-optimized code is used in practice:

- **AWS libcrypto (AWS-LC)**: SLOTHY-optimized X/Ed25519, P256, P384, P521, Keccak, and ML-KEM code has been merged into AWS-LC as part of s2n-bignum.

[1]D. Kostic, H. Becker, J. Harrison, J. Lee, N. Ebeid, and T. Hansen. *Adoption of High-Assurance and Highly Performant Cryptographic Algorithms at AWS.* Real World Crypto 2024. Amazon Web Services, Apr. 2024

# SLOTHY is Used in Practice

SLOTHY-optimized code is used in practice:

- **AWS libcrypto (AWS-LC)**: SLOTHY-optimized X/Ed25519, P256, P384, P521, Keccak, and ML-KEM code has been merged into AWS-LC as part of s2n-bignum.
  - → *"Servicing Trillions of requests a day"*[1]



---

[1]D. Kostic, H. Becker, J. Harrison, J. Lee, N. Ebeid, and T. Hansen. *Adoption of High-Assurance and Highly Performant Cryptographic Algorithms at AWS.* Real World Crypto 2024. Amazon Web Services, Apr. 2024

# SLOTHY is Used in Practice

SLOTHY-optimized code is used in practice:

- **AWS libcrypto (AWS-LC)**: SLOTHY-optimized X/Ed25519, P256, P384, P521, Keccak, and ML-KEM code has been merged into AWS-LC as part of s2n-bignum.
  - → *"Servicing Trillions of requests a day"*[1]
- **s2n-bignum**: AWS' bignum library routinely employs SLOTHY for finding further highly optimized ECC implementations.



---

[1]D. Kostic, H. Becker, J. Harrison, J. Lee, N. Ebeid, and T. Hansen. *Adoption of High-Assurance and Highly Performant Cryptographic Algorithms at AWS.* Real World Crypto 2024. Amazon Web Services, Apr. 2024

SLOTHY-optimized code is used in practice:

- **AWS libcrypto (AWS-LC)**: SLOTHY-optimized X/Ed25519, P256, P384, P521, Keccak, and ML-KEM code has been merged into AWS-LC as part of s2n-bignum.
  - → *"Servicing Trillions of requests a day"*[1]
- **s2n-bignum**: AWS' bignum library routinely employs SLOTHY for finding further highly optimized ECC implementations.
- **Arm EndpointAI**: SLOTHY-optimized code has been deployed to the CMSIS DSP Library for the radix-4 CFFT routines as part of the Arm EndpointAI project



---

[1]D. Kostic, H. Becker, J. Harrison, J. Lee, N. Ebeid, and T. Hansen. *Adoption of High-Assurance and Highly Performant Cryptographic Algorithms at AWS.* Real World Crypto 2024. Amazon Web Services, Apr. 2024

# SLOTHY is Used in Practice

SLOTHY-optimized code is used in practice:

- **AWS libcrypto (AWS-LC)**: SLOTHY-optimized X/Ed25519, P256, P384, P521, Keccak, and ML-KEM code has been merged into AWS-LC as part of s2n-bignum.
  - → *"Servicing Trillions of requests a day"*[1]
- **s2n-bignum**: AWS' bignum library routinely employs SLOTHY for finding further highly optimized ECC implementations.
- **Arm EndpointAI**: SLOTHY-optimized code has been deployed to the CMSIS DSP Library for the radix-4 CFFT routines as part of the Arm EndpointAI project
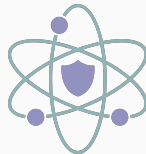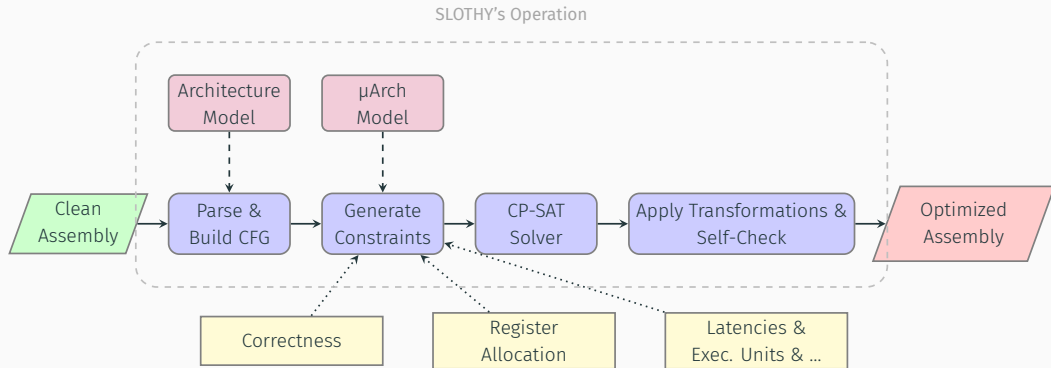- **mlkem-native**: AArch64 assembly routines of ML-KEM are automatically optimized using SLOTHY. **See Matthias talk at OPTIMIST workshop today at 2pm**



---

[1]D. Kostic, H. Becker, J. Harrison, J. Lee, N. Ebeid, and T. Hansen. *Adoption of High-Assurance and Highly Performant Cryptographic Algorithms at AWS.* Real World Crypto 2024. Amazon Web Services, Apr. 2024

# Inside SLOTHY

# Inside SLOTHY: Parsing the Code

## Input Assembly

```
1  // in: src, modulus, const,
        const_twisted
2  mul      dst, src, const
3  sqrdmulh tmp, src, const_twisted
4  mls      dst, tmp, modulus
5  // out: dst
6
```

## Computational Flow Graph

## Computational Flow Graph



- Each instruction has a program position IX.pos (integer variable)
- Naturally, each position can only be assigned once
- For correctness, consumer after producer:
  - I1.pos < I3.pos
  - I2.pos < I3.pos

12

## Computational Flow Graph



**Boolean variables** (reg. is output)

`I1.V0, ..., I1.V31`    `I2.V0, ..., I2.V31`

`I3` needs the same output register as `I1`

**Reg. alloc. constraint**
Exactly one of `I1.V0, ..., I1.V31` is true

**Reg. usage interval** (cond. on boolean variables)

`[I1.pos, I3.pos]`    `[I2.pos, I3.pos]`

**Lifetime constraint**
For each register: Usage intervals must not overlap

13

## Inside SLOTHY: What Else are we Modeling?

- Performance characteristics:
    - Latencies: $\texttt{I1.pos} + 3 < \texttt{I3.pos}$

- Performance characteristics:
  - Latencies: $\mathtt{I1.pos} + 3 < \mathtt{I3.pos}$
  - Occupancy of execution units
  - Throughput, i.e., how long is a unit kept busy?

## Inside SLOTHY: What Else are we Modeling?

- Performance characteristics:
  - Latencies: $\texttt{I1.pos} + 3 < \texttt{I3.pos}$
  - Occupancy of execution units
  - Throughput, i.e., how long is a unit kept busy?
  - Forwarding paths

## Inside SLOTHY: What Else are we Modeling?

- Performance characteristics:
  - Latencies: $\texttt{I1.pos} + 3 < \texttt{I3.pos}$
  - Occupancy of execution units
  - Throughput, i.e., how long is a unit kept busy?
  - Forwarding paths
  - Stalls: In case there is no perfect solution, we model "gaps" in the scheduling
    - **Binary search:** Find solution for a fixed number of stalls; performs (external) binary search to find minimum
    - **Variable size:** Model number of stalls within the constraint model (recommended for small to medium-sized examples)

## Inside SLOTHY: What Else are we Modeling?

- Performance characteristics:
    - Latencies: $\texttt{I1.pos} + 3 < \texttt{I3.pos}$
    - Occupancy of execution units
    - Throughput, i.e., how long is a unit kept busy?
    - Forwarding paths
    - Stalls: In case there is no perfect solution, we model "gaps" in the scheduling
        - **Binary search:** Find solution for a fixed number of stalls; performs (external) binary search to find minimum
        - **Variable size:** Model number of stalls within the constraint model (recommended for small to medium-sized examples)
- Memory dependencies
- Stack spills

## Inside SLOTHY: What Else are we Modeling?

- Performance characteristics:
  - Latencies: $\texttt{I1.pos} + 3 < \texttt{I3.pos}$
  - Occupancy of execution units
  - Throughput, i.e., how long is a unit kept busy?
  - Forwarding paths
  - Stalls: In case there is no perfect solution, we model "gaps" in the scheduling
    - **Binary search:** Find solution for a fixed number of stalls; performs (external) binary search to find minimum
    - **Variable size:** Model number of stalls within the constraint model (recommended for small to medium-sized examples)
- Memory dependencies
- Stack spills

**Note on scalability**: The more properties we model and the more instructions we have, the more complex the constraint problem might be.

$\implies$ At some point, problem becomes infeasible to solve

$\implies$ Workaround: Splitting heuristic optimizes the code piece by piece

# Inside SLOTHY: Software Pipelining



Visualization inspired by https://www.lighterra.com/papers/basicinstructionscheduling/.

## Symbolic Registers in SLOTHY

### Symbolic Registers

- Manual register allocation can be tedious
- SLOTHY will re-allocate registers anyway
- We might as well leave the entire register allocation to SLOTHY

## Symbolic Registers in SLOTHY

### Symbolic Registers

- Manual register allocation can be tedious
- SLOTHY will re-allocate registers anyway
- We might as well leave the entire register allocation to SLOTHY

### Traditional Assembly

```
// Manual register allocation
ldr q0, [x1], #16
ldr q1, [x2], #16
add v0.8h, v0.8h, v1.8h
str q0, [x0], #16
```

### With Symbolic Registers

```
// SLOTHY allocates registers
ldr q<a>, [x1], #16
ldr q<b>, [x2], #16
add v<c>.8h, v<a>.8h, v<b>.8h
str q<c>, [x0], #16
```

## Symbolic Registers in SLOTHY

### Symbolic Registers

- Manual register allocation can be tedious
- SLOTHY will re-allocate registers anyway
- We might as well leave the entire register allocation to SLOTHY

### Traditional Assembly

```
// Manual register allocation
ldr q0, [x1], #16
ldr q1, [x2], #16
add v0.8h, v0.8h, v1.8h
str q0, [x0], #16
```

### With Symbolic Registers

```
// SLOTHY allocates registers
ldr q<a>, [x1], #16
ldr q<b>, [x2], #16
add v<c>.8h, v<a>.8h, v<b>.8h
str q<c>, [x0], #16
```

### Syntax: `type<name>`

- Example: `x<A>` (64-bit), `w<A>` (32-bit view of same register)
- The constraint problem complexity is unchanged (as renaming is done anyway)

## Symbolic Registers in SLOTHY

#### Why Type Prefix is Necessary

- Problem 1: `add A, B, C` - scalar or vector operation?
- Problem 2: Arm has register views, e.g., `x0`–`x30` (64-bit), `w0`–`w30` (32-bit lower half)
  $\implies$ Need to be able to tell apart `add x0, x1, x2` and `add w0, w1, w2`

## Symbolic Registers in SLOTHY

### Why Type Prefix is Necessary

- Problem 1: `add A, B, C` - scalar or vector operation?
- Problem 2: Arm has register views, e.g., `x0`–`x30` (64-bit), `w0`–`w30` (32-bit lower half)
  $\implies$ Need to be able to tell apart `add x0, x1, x2` and `add w0, w1, w2`

### Benefits of using symbolic registers

- Developer does not have to think about register allocation
- Can write code with no valid register allocation without reordering

## Symbolic Registers in SLOTHY

### Why Type Prefix is Necessary
- Problem 1: `add A, B, C` - scalar or vector operation?
- Problem 2: Arm has register views, e.g., `x0–x30` (64-bit), `w0–w30` (32-bit lower half)
  $\implies$ Need to be able to tell apart `add x0, x1, x2` and `add w0, w1, w2`

### Benefits of using symbolic registers
- Developer does not have to think about register allocation
- Can write code with no valid register allocation without reordering

### Disadvantages of using symbolic registers
- Input code is not executable

Caveat: Does not work with split heuristic (large pieces of code)

- Register allocation has to be done globally - if the combined scheduling/allocation problem is too hard you are out of luck
- Workaround: First perform register allocation (disallowing reordering), then perform scheduling using split heuristic

# SLOTHY Basics: Hands-on Assignment 1 + 2

# Installation

## Option 1: Install from PyPI

```
1  python3 -m venv venv
2  source venv/bin/activate
3  pip install slothy
```

# Installation

### Option 1: Install from PyPI

```
1 python3 -m venv venv
2 source venv/bin/activate
3 pip install slothy
```

### Option 2: Development Installation from Source

```
1 git clone https://github.com/slothy-optimizer/slothy.git
2 cd slothy
3 python3 -m venv venv
4 source venv/bin/activate
5 pip install -r requirements.txt
```

### Clone Tutorial Repository
```
$ git clone https://github.com/dop-amin/ches2025-slothy-tutorial
```

### You will find five assignments

- `01basic/` (Basic – Using SLOTHY)
- `02basemul/` (ML-DSA basemul – Using symbolic registers and software pipelining)
- `03keccak/` (Keccak – Optimizing long code with SLOTHY)
- `04instruction/` (EON – Adding instructions to SLOTHY - **requires local SLOTHY clone**)
- `05fusion/` (EOR3 – Using SLOTHY's fusion feature)

Solutions will be presented as a part of this tutorial and can be requested by e-mail.

# How to Use SLOTHY

## Option 1: CLI

```
$ slothy-cli Arm_AArch64 \
    Arm_Cortex_A55 \
    input.s \
    -o output.s \
    -s start_label \
    -e end_label
```

## Option 2: Python Script

```python
1  from slothy import Slothy
2  import slothy.targets.aarch64.aarch64_neon as AArch64_Neon
3  import slothy.targets.aarch64.cortex_a55 as Target_CortexA55
4
5  s = Slothy(AArch64_Neon, Target_CortexA55)
6  s.load_source_from_file("input.s")
7  s.optimize(start="start_label",
8             end="end_label")
9  s.write_source_to_file("output.s")
```

- Vast majority of features available in both
- Today's assignments use Python scripts

## Assignment 1: Polynomial Addition

**Task**
Optimize polynomial addition for ML-KEM using SLOTHY

```
void poly_add(int16_t *r, const int16_t *a, const int16_t *b)
```

- 256 16-bit coefficients
- AArch64 Neon assembly
- Target: Cortex-A55

`01basic/optimize.py`

```python
from slothy import Slothy
import slothy.targets.aarch64.aarch64_neon \
    as AArch64_Neon
import slothy.targets.aarch64.cortex_a55 \
    as Target_CortexA55

def main():
    # TODO: Initialize SLOTHY with AArch64_Neon
    # architecture and Target_CortexA55

    # TODO: Load the source assembly file

    # TODO: Optimize the code between
    # slothy_start and slothy_end markers

    # TODO: Write optimized code to output file

if __name__ == "__main__":
    main()
```

## What You Need to Do

1. Create SLOTHY instance

2. Load `poly_add.s`

3. Call `optimize()` method

4. Save to `poly_add_opt_a55.s`

## Key Files

- `poly_add.s` - input assembly

- `optimize.py` - your script

23

# Input Assembly Structure

```
1   // From: 01basic/poly_add.s
2   mov x3, #4
3   loop:
4     slothy_start:
5     ldr q0, [x1], #128
6     ldr q1, [x1, #-112]
7     // ... (6 more loads from a)
8     ldr q8, [x2], #128
9     ldr q9, [x2, #-112]
10    // ... (6 more loads from b)
11    add v0.8h, v0.8h, v8.8h
12    add v1.8h, v1.8h, v9.8h
13    // ... (6 more adds)
14    str q0, [x0], #128
15    str q1, [x0, #-112]
16    // ... (6 more stores)
17    slothy_end:
18    subs x3, x3, #1
19    b.ne loop
20
```

## Optimization Markers

- Start/end labels (any names work)
- Example: `slothy_start`, `slothy_end`

## Manual Loop Unrolling

- Provides instruction-level parallelism
- Enables effective SLOTHY optimization

## SLOTHY can also optimize loops

- Will be covered in detail later
- Advanced loop features:
  - Software pipelining
  - Automatic loop unrolling

### Basic Import & Setup

```
1  from slothy import Slothy
2
3  # Import architecture
4  import slothy.targets.aarch64.
       aarch64_neon
5      as AArch64_Neon
6
7  # Import target microarchitecture
8  import slothy.targets.aarch64.
       cortex_a55
9      as Target_CortexA55
10
11 # Create SLOTHY instance
12 s = Slothy(AArch64_Neon,
       Target_CortexA55)
```

### Currently Available Targets
AArch64 (Neon):

- cortex_a55
- cortex_a72_frontend
- apple_m1_firestorm_experimental
- apple_m1_icestorm_experimental
- neoverse_n1_experimental
- aarch64_big_experimental

Arm v7-M:

- cortex_m4, cortex_m7

Arm v8.1-M (Helium):

- cortex_m55r1, cortex_m85r1

**`load_source_from_file(filename)`**

- Loads and parses assembly from file

## Core SLOTHY Methods

**`load_source_from_file(filename)`**

- Loads and parses assembly from file

**`write_source_to_file(filename)`**

- Saves optimized assembly with metrics

**`load_source_from_file(filename)`**

- Loads and parses assembly from file

**`write_source_to_file(filename)`**

- Saves optimized assembly with metrics

**`optimize(start=, end=)`**

- `start`/`end` - optimization region markers
- Without markers, SLOTHY optimizes entire file
- Alternative: `optimize_loop(looplabel)` (see next assignment)

Full API documentation: `slothy-optimizer.github.io/slothy/`

## Assignment 2: ML-DSA Basemul

### Task
Use SLOTHY's software pipelining feature to optimize ML-DSA basemul

```
void poly_basemul_montgomery(int16_t *r, const int16_t *a, const int16_t
*b)
```

- Simple loop with symbolic registers
- AArch64 Neon assembly with Montgomery multiplication
- Target: Cortex-A55

```
1   // From: 02basemul/basemul.s
2   modulus          .req v0
3   modulus_twisted  .req v1
4   count            .req x3
5   t0               .req v7
6   .macro montgomery_reduce_long res, inl, inh
7     uzp1   t0.4s, \inl\().4s, \inh\().4s
8     mul    t0.4s, t0.4s, modulus_twisted.4s
9     smlal  \inl\().2d, t0.2s, modulus.2s
10    smlal2 \inh\().2d, t0.4s, modulus.4s
11    uzp2   \res\().4s, \inl\().4s, \inh\().4s
12  .endm
13  .macro pmull dl, dh, a, b
14    smull  \dl\().2d, \a\().2s, \b\().2s
15    smull2 \dh\().2d, \a\().4s, \b\().4s
16  .endm
17  // ... (reg setup, modulus loading)
```

```
1   mov count, #256
2
3   loop_start:
4     ldr q<aa>, [x1], #64
5     ldr q<bb>, [x2], #64
6     pmull v<resl>, v<resh>, v<aa>, v<bb>
7     montgomery_reduce_long v<res>, v<resl>, v<resh>
8     str q<res>, [x0], #64
9
10    ldr q<aa>, [x1, #-48]
11    ldr q<bb>, [x2, #-48]
12    pmull v<resl>, v<resh>, v<aa>, v<bb>
13    montgomery_reduce_long v<res>, v<resl>, v<resh>
14    str q<res>, [x0, #-48]
15
16    // ... (2 more similar iterations)
17    subs count, count, #16
18    cbnz count, loop_start
```

## Optimizing Loops with SLOTHY

**`optimize_loop(loop_lbl)`**

- Optimizes a loop starting at a given label
- Automatically detects loop structure
- Example: `slothy.optimize_loop("loop_start")`

## Optimizing Loops with SLOTHY

### `optimize_loop(loop_lbl)`

- Optimizes a loop starting at a given label
- Automatically detects loop structure
- Example: `slothy.optimize_loop("loop_start")`

### Supported Loop Patterns (AArch64)

- Counter decrement: `sub[s] <reg>, <reg>, #<imm>`
- Followed by conditional branch:
    - `cbnz/cbz <reg>, <loop_lbl>` (compare and branch)
    - `b.<cond> <loop_lbl>` (any Arm condition: ne, eq, ge, lt, etc.)

# Optimizing Loops with SLOTHY

## `optimize_loop(loop_lbl)`

- Optimizes a loop starting at a given label
- Automatically detects loop structure
- Example: `slothy.optimize_loop("loop_start")`

## Supported Loop Patterns (AArch64)

- Counter decrement: `sub[s] <reg>, <reg>, #<imm>`
- Followed by conditional branch:
    - `cbnz/cbz <reg>, <loop_lbl>` (compare and branch)
    - `b.<cond> <loop_lbl>` (any Arm condition: ne, eq, ge, lt, etc.)
- **Note:** Loop patterns are extensible - additional patterns can be added to the architecture model

# Optimizing Loops with SLOTHY

## `optimize_loop(loop_lbl)`

- Optimizes a loop starting at a given label
- Automatically detects loop structure
- Example: `slothy.optimize_loop("loop_start")`

## Supported Loop Patterns (AArch64)

- Counter decrement: `sub[s] <reg>, <reg>, #<imm>`
- Followed by conditional branch:
  - `cbnz/cbz <reg>, <loop_lbl>` (compare and branch)
  - `b.<cond> <loop_lbl>` (any Arm condition: ne, eq, ge, lt, etc.)
- **Note:** Loop patterns are extensible - additional patterns can be added to the architecture model

## Limitations

- Loop counter must be decremented by a constant
- Loop body must be straight-line code (no branches inside)
- Subtraction must immediately precede the branch instruction

# SLOTHY Configuration Options

### Setting Configuration Options
Configuration options are set after creating the SLOTHY instance:

### Example: Setting Configuration Options

```
1  # Create SLOTHY instance
2  slothy = Slothy(Architecture, Target)
3
4  # Configure SLOTHY settings
5  slothy.config.<option1> = <value1>
6  slothy.config.<option2> = <value2>
7  slothy.config.<option3>.<suboption> = <value3>
8
9  # ... load, optimize, write
10
```

Full configuration reference:

`slothy-optimizer.github.io/slothy/apidocs/slothy/slothy.core.config.html`

### slothy.config.sw_pipelining.enabled

Enable software pipelining optimization for loop code blocks.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Essential for optimizing loops with overlapping iterations

### slothy.config.sw_pipelining.allow_pre

Allow 'early' instructions to be pulled forward from future iterations.

- **Default:** `True`
- **Type:** Boolean (`True`|`False`)
- **Usage:** Enables forward movement of instructions from iteration N+1 to N

### slothy.config.sw_pipelining.allow_post

Allow 'late' instructions to be deferred to later iterations.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Enables backward movement of instructions from iteration N to N+1

### slothy.config.sw_pipelining.unroll

The unrolling factor to use for software pipelining optimization.

- Default: 1
- Type: Integer (positive number)
- Usage: Higher values enable more aggressive pipelining optimizations

### slothy.config.sw_pipelining.optimize_preamble

Perform a separate optimization pass for the loop preamble section.

- Default: `True`
- Type: Boolean (`True|False`)
- Usage: Optimizes setup instructions before the main pipelined loop

### slothy.config.sw_pipelining.optimize_postamble

Perform a separate optimization pass for the loop postamble section.

- **Default:** `True`
- **Type:** Boolean (`True|False`)
- **Usage:** Optimizes cleanup instructions after the main pipelined loop

## Performance Configuration

### `slothy.config.variable_size`

Model the number of stalls as a parameter in the constraint satisfaction problem.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Enable for small/medium code where solver can minimize stalls directly; if `False`: External binary search finds minimum stalls

### `slothy.config.constraints.stalls_first_attempt`

The initial number of stalls to attempt during optimization.

- Default: `0`
- Type: Integer (non-negative number)
- Usage: Set higher when minimum stall count is known from experience

## Input/Output Configuration

### slothy.config.outputs

List of architectural registers that must be preserved as function outputs.

- Default: `[]`
- Type: List of strings (register names or `flags`)
- Usage: Critical for correctness - prevents clobbering output registers

### slothy.config.inputs_are_outputs

Treat all input registers as outputs that must be preserved.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Typically used for loop optimization

### slothy.config.reserved_regs

Set of architectural registers excluded from register renaming and allocation.

- **Default:** `["sp"]` (AArch64)
- **Type:** List of strings (register names)
- **Note:** Overwrites default reserved registers for the target architecture

**Note:** Reserved registers are "locked" by default - they won't be introduced during renaming, but existing uses won't be touched either.

[30 minutes hands-on exercise (Assignment 1 and 2)]
See the `README.md` for hints



**Tutorial Assignments**
github.com/dop-amin/ches2025-slothy-tutorial



**Tutorial Slides**
kannwischer.eu/talks/20250914_slothy.pdf

`01basic/optimize.py`

```python
1  from slothy import Slothy
2  import slothy.targets.aarch64.aarch64_neon as AArch64_Neon
3  import slothy.targets.aarch64.cortex_a55 as Target_CortexA55
4
5  def main():
6      # Initialize SLOTHY
7      slothy = Slothy(AArch64_Neon, Target_CortexA55)
8
9      # Load the source assembly file
10     slothy.load_source_from_file("poly_add.s")
11
12     # Optimize the code between markers
13     slothy.optimize(start="slothy_start", end="slothy_end")
14
15     # Write optimized code to output file
16     slothy.write_source_to_file("poly_add_opt_a55.s")
17
18  if __name__ == "__main__":
19      main()
```

## Assignment 1: Optimized Result

```
1    slothy_start:
2                                        // Instructions:     32
3                                        // Expected cycles: 49
4                                        // Expected IPC:     0.65
5                                        //
6                                        // Cycle bound:      49.0
7                                        // IPC bound:        0.65
8                                        //
9                                        // Wall time:        0.07s
10                                       // User time:        0.07s
11                                       //
12                                       // -------------- cycle (expected) --------------->
13                                       // 0                        25
14                                       // |-----------------------|-----------------------
15       ldr q27, [x1], #128             // *................................................
16       ldr q8, [x2], #128              // ..*..............................................
17       ldr q30, [x1, #-112]            // ....*............................................
18       add v9.8H, v27.8H, v8.8H        // ......*..........................................
19       ldr q23, [x2, #-112]            // .......*.........................................
20       ldr q1, [x2, #-96]              // .........*.......................................
21       add v4.8H, v30.8H, v23.8H       // ...........*.....................................
22       str q9, [x0], #128              // ............*....................................
23       ldr q15, [x1, #-96]             // .............*...................................
24       ldr q10, [x1, #-80]             // ..............*..................................
25       ldr q7, [x2, #-80]              // ...............*.................................
26       add v0.8H, v15.8H, v1.8H        // .................*...............................
27       // ... (continues)
28   slothy_end:
29
```

42

02basemul/optimize.py

```
1   from slothy import Slothy
2   import slothy.targets.aarch64.aarch64_neon as AArch64_Neon
3   import slothy.targets.aarch64.cortex_a55 as Target_CortexA55
4
5   def main():
6       # Initialize SLOTHY
7       slothy = Slothy(AArch64_Neon, Target_CortexA55)
8
9       # Enable software pipelining
10      slothy.config.sw_pipelining.enabled = True
11      slothy.config.variable_size = True
12      slothy.config.constraints.stalls_first_attempt = 32
13
14      # Load source and optimize loop
15      slothy.load_source_from_file("basemul.s")
16      slothy.optimize_loop("loop_start")
17      slothy.write_source_to_file("basemul_opt_a55.s")
18
19  if __name__ == "__main__":
20      main()
```

```
1           ldr q7, [x2, #32]        // *.............................
2           sub count, count, #16
3    loop_start:
4                                    // Instructions:    40
5                                    // Expected cycles: 48
6                                    // Expected IPC:    0.83
7                                    //
8                                    // Cycle bound:      21.0
9                                    // IPC bound:        1.90
10                                   //
11                                   // Wall time:      5.36s
12                                   // User time:      5.36s
13                                   //
14                                   // ------------- cycle (expected) --------------->
15                                   // 0                            25
16                                   // |-----------------------|---------------------
17          ldr q23, [x2], #64       // *.................................................
18          ldr q10, [x1], #64       // ..*...............................................
19          ldr q4, [x1, #-32]       // ....*.............................................
20          smull v2.2D, v10.2S, v23.2S   // ......*.......................................
21          smull2 v13.2D, v10.4S, v23.4S // .......*......................................
22          smull v5.2D, v4.2S, v7.2S     // ........*.....................................
23          smull2 v14.2D, v4.4S, v7.4S   // .........*....................................
24          ldr q4, [x1, #-16]       // ..........*.......................................
25          uzp1 v31.4S, v2.4S, v13.4S    // ...........*..................................
26          ldr q24, [x2, #-16]      // ............*.....................................
27          // ... (continues)
28
```

# Did our code actually get faster?

| Assignment | Before (cycles) | After (cycles) | Speedup |
|---|---|---|---|
| 1) Poly Add (A55) | 228 | 204 | $1.12\times$ |
| 2) Poly Basemul (A55) | 1512 | 822 | $1.84\times$ |

Table 1: Performance on Cortex-A55 cores

# Heuristics and Register Spilling

Hardness of the problems SLOTHY deals with depends on a number of factors:

- Number of instructions in the assembly
- Complexity of $\mu$Arch model
- Is software pipelining/spilling wanted?

## Splitting Heuristic

Usually: SLOTHY considers the input code as one, large chunk. Optimizing this chunk may not terminate if there are too many instructions (i.e., the problem gets too hard).

# Splitting Heuristic

**Usually**: SLOTHY considers the input code as one, large chunk. Optimizing this chunk may not terminate if there are too many instructions (i.e., the problem gets too hard).

**Idea**: Pick a window containing a modest number of instructions and slide this window across the code, then repeat.
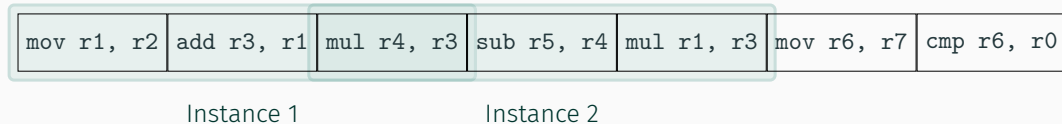
| mov r1, r2 | add r3, r1 | mul r4, r3 | sub r5, r4 | mul r1, r3 | mov r6, r7 | cmp r6, r0 |

       Instance 1

Usually: SLOTHY considers the input code as one, large chunk. Optimizing this chunk may not terminate if there are too many instructions (i.e., the problem gets too hard).

Idea: Pick a window containing a modest number of instructions and slide this window across the code, then repeat.

| mov r1, r2 | add r3, r1 | mul r4, r3 | sub r5, r4 | mul r1, r3 | mov r6, r7 | cmp r6, r0 |
|---|---|---|---|---|---|---|

Instance 1          Instance 2

# Splitting Heuristic

Usually: SLOTHY considers the input code as one, large chunk. Optimizing this chunk may not terminate if there are too many instructions (i.e., the problem gets too hard).
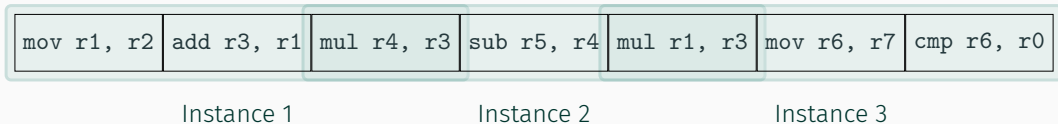
Idea: Pick a window containing a modest number of instructions and slide this window across the code, then repeat.

| mov r1, r2 | add r3, r1 | mul r4, r3 | sub r5, r4 | mul r1, r3 | mov r6, r7 | cmp r6, r0 |
|---|---|---|---|---|---|---|

       Instance 1             Instance 2             Instance 3

**Usually**: SLOTHY considers the input code as one, large chunk. Optimizing this chunk may not terminate if there are too many instructions (i.e., the problem gets too hard).
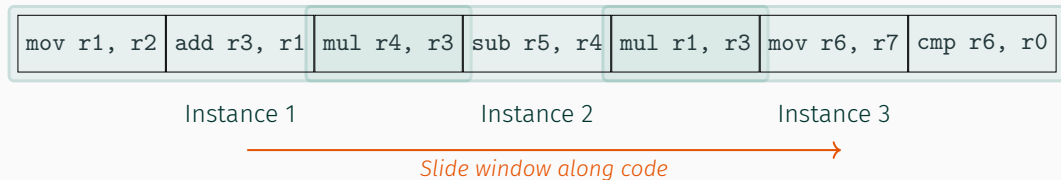
**Idea**: Pick a window containing a modest number of instructions and slide this window across the code, then repeat.



```
mov r1, r2 | add r3, r1 | mul r4, r3 | sub r5, r4 | mul r1, r3 | mov r6, r7 | cmp r6, r0
```

Instance 1          Instance 2          Instance 3

*Slide window along code*

### slothy.config.split_heuristic

Trade-off between runtime and optimality by splitting code blocks into subchunks.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Enable when optimization of large code blocks fails

## Splitting Heuristic: Useful Options ii

### slothy.config.split_heuristic_factor

Number of subchunks to split each code block into for optimization.

- **Default:** 2
- **Type:** Integer (positive number)
- **Usage:** Only meaningful when `split_heuristic=True`

### slothy.config.split_heuristic_stepsize

Increment for the sliding window as fraction of total code size.

- Default: `None`
- Type: Float (0.0 to 1.0)
- Usage: Controls overlap between optimization windows

### `slothy.config.split_heuristic_repeat`

Number of times the splitting heuristic procedure should be repeated.

- **Default:** 1
- **Type:** Integer (positive number)
- **Usage:** Facilitate more interleaving; improves output code

# Interleaving Heuristic i

Whenever the code-paths to be interleaved are very far apart, establishing a coarse interleaving on the CFG can be helpful to facilitate later optimization.

## slothy.config.split_heuristic_preprocess_naive_interleaving

Interleave instructions by lowest depth without applying register renaming.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Helps when code paths to be interleaved are far apart

**slothy.config.**
**split_heuristic_preprocess_naive_interleaving_strategy**

Strategy for naive interleaving preprocessing of distant code paths.

- Default: `"depth"`
- Type: String (`"depth"`|`"alternate"`)
- Usage: Choose interleaving pattern for coarse CFG optimization

# Spill-code Generation

Problem: When writing a symbolic implementation, register allocation is unclear.

# Spill-code Generation

Problem: When writing a symbolic implementation, register allocation is unclear.

- How many registers will be required?
- Are there enough registers?
- What are we going to do when we run out of registers?

# Spill-code generation

Solution: Let SLOTHY take care of this. SLOTHY can ...

- ...introduce spills/restores whenever it runs out of registers.

## Spill-code generation

Solution: Let SLOTHY take care of this. SLOTHY can ...

- · ...introduce spills/restores whenever it runs out of registers.
- · ...minimize the number of spills it introduces as an optimization target.

## Spill-code generation

Solution: Let SLOTHY take care of this. SLOTHY can …

- …introduce spills/restores whenever it runs out of registers.
- …minimize the number of spills it introduces as an optimization target.
- …absorb superfluous spills in existing, non-symbolic code in case a register allocation without these spills exists.

# Spill-code generation

Solution: Let SLOTHY take care of this. SLOTHY can …

- …introduce spills/restores whenever it runs out of registers.
- …minimize the number of spills it introduces as an optimization target.
- …absorb superfluous spills in existing, non-symbolic code in case a register allocation without these spills exists.
- …*not* minimize spills AND reorder instructions simultaneously → State explosion

### slothy.config.allow_spills

Allow SLOTHY to introduce stack spills when register pressure is too high.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Essential for symbolic assembly with high register pressure

### slothy.config.minimize_spills

Minimize the number of stack spills as the solver's optimization objective.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Enable when spill count minimization is the primary goal

### slothy.config.absorb_spills

Remove existing spills from non-symbolic code when better register allocation exists. Does *not* work with splitting heuristic.

- Default: `True`
- Type: Boolean (`True|False`)
- Usage: Cleanup existing assembly with unnecessary spills

## slothy.config.constraints.functional_only

Disable modeling of latencies and functional units.

- Default: `False`
- Type: Boolean (`True`|`False`)
- Usage: Reduce complexity of the optimization problem when, e.g., focusing on register allocation

**slothy.config.constraints.allow_reordering**

Allow reordering of instructions.

- **Default:** `True`
- **Type:** Boolean (`True|False`)
- **Usage:** Reduce complexity of the optimization problem when, e.g., focusing on register allocation

```
1  ldr X<count>, [sp, #STACK_OFFSET_COUNT] //
        @slothy:reads=STACK_OFFSET_COUNT
2  ldr X<const>, [X<caddr>, X<count>, LSL #3]
3  add X<count>, X<count>, #1
4  cmp X<count>, #(KECCAK_F1600_ROUNDS-1)
5  str X<count>, [sp, #STACK_OFFSET_COUNT] //
        @slothy:writes=STACK_OFFSET_COUNT
6
```

- Reads/writes marked through `@slothy` tags
- Enables safely re-scheduling interdependent reads/writes
- Internally: "hint"-registers

## Assignment 3: Keccak Permutation

### Task
Use SLOTHY's spill-code generation & splitting heuristic to optimize the Keccak permutation.

```
void keccak_f1600_x4_v8a_hybrid_slothy_symbolic(void* states)
```

- Symbolic implementation; not enough registers
- Large number of instructions
- AArch64 Scalar+Neon hybrid assembly
- Target: Cortex-A55

## Assignment 3: Input Assembly

```
 1  // Many macros...
 2  keccak_f1600_x4_v8a_hybrid_slothy_symbolic:
 3      // Some code
 4   initial:
 5      scalar_round_initial    // @slothy:interleaving_class=0
 6      scalar_round_noninitial // @slothy:interleaving_class=0
 7      vector_round            // @slothy:interleaving_class=1
 8   loop:
 9      scalar_round_noninitial // @slothy:interleaving_class=0
10      scalar_round_noninitial // @slothy:interleaving_class=0
11      vector_round            // @slothy:interleaving_class=1
12   loop_end:
13      ble loop
14      // More code
15      b initial
16   done:
17      // More code
18      ret
19
```

- Two intervals to optimize
- // @slothy:... tags
  - → Eased coarse interleaving
  - → Memory dependencies

# Assignment 3: Python Template

```python
from slothy import Slothy
import slothy.targets.aarch64.aarch64_neon as AArch64_Neon
import slothy.targets.aarch64.cortex_a55 as Target_CortexA55

def main():
    slothy = Slothy(AArch64_Neon, Target_CortexA55)
    # Load the source assembly file
    slothy.load_source_from_file("keccak.s")
    # ...
    # --- Register Allocation
    # TODO: complete
    slothy.write_source_to_file("keccak_alloc_a55.s")
    # ...
    # --- Optimize
    slothy.load_source_from_file("keccak_alloc_a55.s")
    # TODO: complete
    # Write optimized code to output file
    slothy.write_source_to_file("keccak_opt_a55.s")
```

Two separate goals for optimization

- Register allocation
- Instruction Scheduling

[30 minutes hands-on exercise (Assignment 3)]
See the `README.md` for hints



Tutorial Assignments
github.com/dop-amin/ches2025-slothy-tutorial



Tutorial Slides
kannwischer.eu/talks/20250914_slothy.pdf

65

# Assignment 3: Solution, Setup

03keccak/optimize.py

```python
1  from slothy import Slothy
2  import slothy.targets.aarch64.aarch64_neon as AArch64_Neon
3  import slothy.targets.aarch64.cortex_a55 as Target_CortexA55
4
5  def main():
6      # Initialize SLOTHY
7      slothy = Slothy(AArch64_Neon, Target_CortexA55)
8      # Load the source assembly file
9      slothy.load_source_from_file("keccak.s")
10     # Common config
11     slothy.config.selftest = False
12     slothy.config.timeout = 180
13     slothy.config.constraints.stalls_first_attempt = 32
14     slothy.config.reserved_regs = ["sp"]
15     slothy.config.outputs = ["hint_STACK_OFFSET_COUNT"]  # preserve count
16     slothy.config.inputs_are_outputs = True
17     slothy.config.variable_size = True
18     slothy.config.with_preprocessor = True
19     common_conf = slothy.config.copy()
```

# Assignment 3: Solution, Register Allocation

03keccak/optimize.py

```
1     slothy.config.constraints.functional_only = True
2     slothy.config.constraints.allow_reordering = False
3     slothy.config.constraints.allow_spills = True
4     slothy.config.constraints.minimize_spills = True
5
6     # Call to the optimizer
7     slothy.optimize(start="loop", end="loop_end")
8     slothy.optimize(start="initial", end="loop")
9
10    slothy.write_source_to_file("keccak_alloc_a55.s")
11
12    # Reset configuration to common
13    slothy.config = common_conf.copy()
```

03keccak/optimize.py

```
1    slothy.load_source_from_file("keccak_alloc_a55.s")
2
3    # Configure optimization parameters
4    slothy.config.split_heuristic = True
5    slothy.config.split_heuristic_preprocess_naive_interleaving = True
6    slothy.config.split_heuristic_preprocess_naive_interleaving_strategy = "alternate"
7
8    slothy.config.split_heuristic_factor = 12
9    slothy.config.split_heuristic_repeat = 2
10   slothy.config.split_heuristic_stepsize = 0.05
11   slothy.config.absorb_spills = False   # Does not work with splitting
12
13   # Call to the optimizer
14   slothy.optimize(start="loop", end="loop_end")
15   slothy.optimize(start="initial", end="loop")
16
17   # Write optimized code to output file
18   slothy.write_source_to_file("keccak_opt_a55.s")
```

| Assignment | Before | | After | | |
|---|---|---|---|---|---|
| | Cycles | IPC | Cycles | IPC | Speedup |
| Keccak | 8212 | 1.18 | 5385 | 1.81 | $1.53\times$ |

Table 2: Keccak performance on Cortex-A55 cores

# Modelling Architectures & Microarchitectures

## Extending SLOTHY: Architecture & Microarchitecture Models

### Motivation

- SLOTHY is built to be extensible
- We would love to see new architectures or extensions to existing ones

## Extending SLOTHY: Architecture & Microarchitecture Models

### Motivation

- SLOTHY is built to be extensible
- We would love to see new architectures or extensions to existing ones

### When You'll Need This Knowledge

- Adding support for a **new architecture** (e.g., RISC-V)
- Adding a **new microarchitecture** (e.g., Arm Cortex-A510)
- Adding **missing instructions** to existing architecture or microarchitecture models
- **Fixing or refining** performance characteristics

## Extending SLOTHY: Architecture & Microarchitecture Models

### Motivation
- SLOTHY is built to be extensible
- We would love to see new architectures or extensions to existing ones

### When You'll Need This Knowledge
- Adding support for a **new architecture** (e.g., RISC-V)
- Adding a **new microarchitecture** (e.g., Arm Cortex-A510)
- Adding **missing instructions** to existing architecture or microarchitecture models
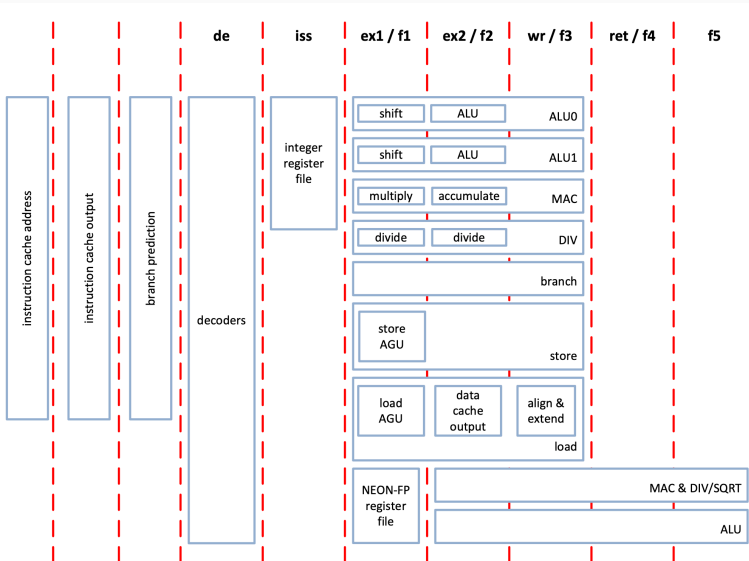- **Fixing or refining** performance characteristics

### The Reality of SLOTHY Models
- SLOTHY's architecture and microarchitecture models are **lazily built**
- Only commonly used instructions are modeled
- Models grow as users need them (AArch64 is the most mature)
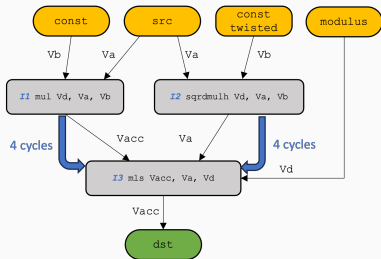
# Understanding the Microarchitecture

Starting Point: Software Optimization Guides (SWOG)

- Manufacturers often provide optimization guides
- Contains pipeline diagrams, latency tables, and throughput information
- Example: Cortex-A55 `https://developer.arm.com/documentation/EPM128372/latest/`

# Cortex-A55 Pipeline

| Instruction group | AArch64 instructions | Exec latency | Execution throughput |
|---|---|---|---|
| ASIMD multiply | MUL, SQDMULH, SQRDMULH | 4 | 2* |

\* If the instruction has Q-form, the Q-form of the instruction can only be dual issued as

instruction 0 and execution throughput is 1.

- Q-form: `mul v0.4s, v1.4s, v2.4s` (128-bit operation)
- D-form: `mul v0.2s, v1.2s, v2.2s` (64-bit operation)

### How constraint models are built from the microarchitecture model
Note: You don't need to implement this yourself - SLOTHY handles it internally.

Each instruction `I` is assigned:

- **Functional unit(s):** `Unit(I)` - where it executes
- **Usage time:** `block(I)` - how long each unit is occupied (inverse throughput)
- **Latency:** `lat(I)` - when results are available

# Constraint Model Intuition

## How constraint models are built from the microarchitecture model
**Note:** You don't need to implement this yourself - SLOTHY handles it internally.

Each instruction `I` is assigned:

- **Functional unit(s):** `Unit(I)` - where it executes
- **Usage time:** `block(I)` - how long each unit is occupied (inverse throughput)
- **Latency:** `lat(I)` - when results are available

## Latency $\neq$ Usage Time

- Latency can be **smaller** or **larger** than usage time

- Example: `VADD` on Cortex-M55:

  - Occupies vector unit for 2 cycles (usage time)
  - Result available after 1 cycle (latency)

## Core Microarchitectural Constraints

1. Latency Constraints
   - Consumer must wait for producer's results
   - `consumer.pos` $\geq$ `producer.pos + latency`

## Core Microarchitectural Constraints

### 1. Latency Constraints

- Consumer must wait for producer's results
- `consumer.pos` $\geq$ `producer.pos + latency`

### 2. Execution Unit + Throughput Constraints

- Instructions cannot overlap on the same unit
- For each unit maintain a list of usage intervals: `[pos, pos + inverse_throughput]`

# Core Microarchitectural Constraints

1. Latency Constraints
   - Consumer must wait for producer's results
   - `consumer.pos` $\geq$ `producer.pos + latency`

2. Execution Unit + Throughput Constraints
   - Instructions cannot overlap on the same unit
   - For each unit maintain a list of usage intervals: `[pos, pos + inverse_throughput]`

**Warning: SWOG and SLOTHY express throughput differently!**
SWOG shows overall throughput, we need inverse throughput per unit!
Commonly:

   - SWOG: "Throughput = 8" $\Rightarrow$ usually 8 units $\Rightarrow$ SLOTHY: inverse throughput = 1
   - SWOG: "Throughput = 1/2" $\Rightarrow$ 1 unit, busy for 2 cycles $\Rightarrow$ SLOTHY: inverse throughput = 2

$\Longrightarrow$ $\text{tp}_{\text{SLOTHY}} = \frac{\#\text{units}}{\text{tp}_{\text{SWOG}}}$

## Modeling Q-form vs D-form Throughput

### Special case of the Arm Cortex-A55

- D-form (64-bit): throughput = 2 (can dual-issue)
- Q-form (128-bit): throughput = 1 (single-issue only)

## Modeling Q-form vs D-form Throughput

### Special case of the Arm Cortex-A55

- D-form (64-bit): throughput = 2 (can dual-issue)
- Q-form (128-bit): throughput = 1 (single-issue only)

### SLOTHY's Solution: Model the vector unit as two virtual units: VEC0 and VEC1

- **D-form instructions:** Occupy either VEC0 **or** VEC1
- **Q-form instructions:** Occupy **both** VEC0 **and** VEC1

## Modeling Q-form vs D-form Throughput

### Special case of the Arm Cortex-A55

- D-form (64-bit): throughput = 2 (can dual-issue)
- Q-form (128-bit): throughput = 1 (single-issue only)

### SLOTHY's Solution: Model the vector unit as two virtual units: VEC0 and VEC1

- **D-form instructions:** Occupy either VEC0 or VEC1
- **Q-form instructions:** Occupy both VEC0 and VEC1

#### Keep in mind: SLOTHY models are approximate!

- We do not try to model each aspect of the microarchitecture.
- **Goal**: Simple enough model that achieves good performance!
- Above's example: Currently SLOTHY does not model the D-form of multiplications as we did not need it so far.

# Example Walkthrough: Neon Multiplication

### What We'll Cover
Walk through how the Neon `mul` instruction is modeled in SLOTHY:

1. Architecture model (instruction syntax and semantics)

2. Microarchitecture model (performance characteristics)

3. How they work together

`slothy/targets/aarch64/aarch64_neon.py`

```python
class vmul(AArch64Instruction):
    pattern = "mul <Vd>.<dt0>, <Va>.<dt1>, <Vb>.<dt2>"
    inputs = ["Va", "Vb"]
    outputs = ["Vd"]

```

What This Defines

- **Pattern:** How to parse/emit the instruction
- **Inputs:** Source operands (Va, Vb)
- **Outputs:** Destination operand (Vd)
- **Placeholders:** `<dt0>, <dt1>, <dt2>` for data types (e.g., 4s, 8h)
- Parsing complexity is hidden away in `AArch64Instruction` class

# Microarchitecture Model Structure

`slothy/targets/aarch64/cortex_a55.py` - Basic Structure

```python
1    class ExecutionUnit(Enum):
2        SCALAR_ALU0 = 1
3        SCALAR_ALU1 = 2
4        SCALAR_MAC = 3
5        SCALAR_LOAD = 4
6        SCALAR_STORE = 5
7        VEC0 = 6
8        VEC1 = 7
9
10   execution_units = {
11       # Map instruction classes to execution units
12       # ...
13   }
14
15   inverse_throughput = {
16       # Map instruction classes to inverse throughput
17       # ...
18   }
19
20   default_latencies = {
21       # Map instruction classes to result latency
22       # ...
23   }
24
```

## Adding Performance Characteristics

```
1    # From slothy/targets/aarch64/cortex_a55.py
2    execution_units = {
3        # ...
4        vmul: [[ExecutionUnit.VEC0, ExecutionUnit.VEC1]], # Q-form uses both VEC0 and VEC1
5    }
6
7    inverse_throughput = {
8        # ...
9        vmul: 1, # Occupies units for 1 cycle
10   }
11
12   default_latencies = {
13       # ...
14       vmul: 4, # Result available after 4 cycles
15   }
```

## Alternative: Modeling Q-form vs D-form

```
1    # More precise (not currently in model):
2    is_qform_form_of(vmul): [[ExecutionUnit.VEC0, ExecutionUnit.VEC1]],  # VEC0 AND VEC1
3    is_dform_form_of(vmul): [ExecutionUnit.VEC0, ExecutionUnit.VEC1],  # VEC0 OR VEC1
```

# Grouping Instructions with Class Hierarchies

## Architecture Model: Using Inheritance for Groups

```python
# slothy/targets/aarch64/aarch64_neon.py
class AArch64NeonLogical(AArch64Instruction):
    pass

class veor(AArch64NeonLogical):
    pattern = "eor <Vd>.<dt0>, <Va>.<dt1>, <Vb>.<dt2>"
    inputs = ["Va", "Vb"]
    outputs = ["Vd"]

class vbic(AArch64NeonLogical):
    pattern = "bic <Vd>.<dt0>, <Va>.<dt1>, <Vb>.<dt2>"
    inputs = ["Va", "Vb"]
    outputs = ["Vd"]
```

# Grouping Instructions with Class Hierarchies

## Architecture Model: Using Inheritance for Groups

```python
1  # slothy/targets/aarch64/aarch64_neon.py
2  class AArch64NeonLogical(AArch64Instruction):
3      pass
4
5  class veor(AArch64NeonLogical):
6      pattern = "eor <Vd>.<dt0>, <Va>.<dt1>, <Vb>.<dt2>"
7      inputs = ["Va", "Vb"]
8      outputs = ["Vd"]
9
10 class vbic(AArch64NeonLogical):
11     pattern = "bic <Vd>.<dt0>, <Va>.<dt1>, <Vb>.<dt2>"
12     inputs = ["Va", "Vb"]
13     outputs = ["Vd"]
```

## Microarchitectural Model: Assign Properties to Groups

```python
1  # slothy/targets/aarch64/cortex_a55.py
2  default_latencies = {
3      # ...
4      AArch64NeonLogical: 1
5  }
```

# Microarchitectural Model: Forwarding Paths

### Fast Result Forwarding
Some CPUs have special paths that allow results to be used faster in specific cases

### Example: Cortex-A72 mla→mla Forwarding

```python
1   # slothy/targets/aarch64/cortex_a72_frontend.py
2   def get_latency(src, out_idx, dst):
3       # Default latency (e.g., vmla has latency 5)
4       latency = lookup_multidict(default_latencies, src)
5
6       # forwarding paths and other special classes
7       instclass_src = find_class(src)
8       instclass_dst = find_class(dst)
9
10      # Fast mla->mla forwarding: reduces latency to 1
11      if (instclass_src == vmla and instclass_dst == vmla
12          and src.args_in_out[0] == dst.args_in_out[0]): # Same accumulator
13          return 1  # Instead of 5!
14
15      return latency
16
```

# Microarchitectural Model: Slot Constraints

### Issue Restrictions
Some instructions can only issue in specific issue slots

### Example: Cortex-A55 Q-form Restrictions

```python
1   # slothy/targets/aarch64/cortex_a55.py
2   def add_further_constraints(slothy):
3       # fcsel
4       slothy.restrict_slots_for_instructions_by_class(
5           [fcsel_dform], [0]
6       )
7
8       # Q-form vector instructions on slot 0 only
9       slothy.restrict_slots_for_instructions_by_property(
10          Instruction.is_q_form_vector_instruction, [0]
11      )
12
```

# Instruction Fusion & Splitting

# Instruction Fusion & Splitting

Two complementary optimization techniques:

# Instruction Fusion & Splitting

Two complementary optimization techniques:

- **Instruction Fusion**: Combine multiple "simple" instructions into one complex instruction

## Instruction Fusion & Splitting

Two complementary optimization techniques:

- **Instruction Fusion**: Combine multiple "simple" instructions into one complex instruction
- **Instruction Splitting**: Break "complex" instructions into simpler ones for better scheduling

Two complementary optimization techniques:

- **Instruction Fusion**: Combine multiple "simple" instructions into one complex instruction
- **Instruction Splitting**: Break "complex" instructions into simpler ones for better scheduling

### Idea: Allow certain simple instruction replacements

Different microarchitectures have different capabilities. Let SLOTHY adapt code to better utilize target-specific features!

# Example: Instruction Splitting on Cortex-M7 (dual issue)

Without Splitting (12 cycles):

```
1  ldm    r0, {r2-r9}   // *...........
2  uadd16 r8, r8, r1    // ....*.......
3  uadd16 r5, r5, r1    // .....*......
4  uadd16 r7, r7, r1    // ......*.....
5  uadd16 r6, r6, r1    // .......*....
6  uadd16 r4, r4, r1    // ........*...
7  uadd16 r3, r3, r1    // .........*..
8  uadd16 r2, r2, r1    // ..........*.
9  uadd16 r9, r9, r1    // ...........*
```

# Example: Instruction Splitting on Cortex-M7 (dual issue)

**Without Splitting (12 cycles):**

```
1  ldm r0, {r2-r9}     // *...........
2  uadd16 r8, r8, r1   // ....*.......
3  uadd16 r5, r5, r1   // .....*......
4  uadd16 r7, r7, r1   // ......*.....
5  uadd16 r6, r6, r1   // ........*....
6  uadd16 r4, r4, r1   // .........*...
7  uadd16 r3, r3, r1   // ..........*..
8  uadd16 r2, r2, r1   // ...........*.
9  uadd16 r9, r9, r1   // ............*
```

**With Splitting (9 cycles):**

```
1  ldr r10, [r0, #0]     // *........
2  uadd16 r2, r10, r1    // .*.......
3  ldr r11, [r0, #28]    // .*.......
4  uadd16 r9, r11, r1    // ..*......
5  ldr r14, [r0, #4]     // ..*......
6  uadd16 r3, r14, r1    // ...*.....
7  ldr r4, [r0, #8]      // ...*.....
8  // ... (dual-issued)
```

### 📈 Performance Impact

**33% speedup** by enabling dual-issue of loads and arithmetic operations!

# Splitting Callback

```
1   def ldm_interval_splitting_cb():
2     def core(inst, t, log=None):
3       ptr = inst.args_in[0]
4       regs = inst.args_out
5       width = inst.width
6       t.inst = []
7       offset = 0
8       for r in regs:
9         ldr = Armv7mInstruction.build(
10          ldr_with_imm, {"width": width, "Rd": r, "Ra": ptr, "imm": f"#{off}"})
11        ldr.pre_index = offset
12        t.inst.append(ldr)
13        offset += 4
14        ldr_src = (SourceLine(ldr.write()).add_tags(inst.source_line.tags).add_comments(inst.source_line.comments))
15        ldr.source_line = ldr_src
16      t.changed = True
17      return True
18    return core
```

## Splitting & Fusion: Limitations

- Splitting/Fusion can either be enabled/disabled; SLOTHY is not able to determine dynamically if applying the heuristic would be beneficial
  - → State explosion

## Splitting & Fusion: Limitations

- Splitting/Fusion can either be enabled/disabled; SLOTHY is not able to determine dynamically if applying the heuristic would be beneficial
  - → State explosion
- Splitting/Fusion is not guaranteed to succeed, e.g., in case of the transformation requiring more registers

## Splitting & Fusion: Limitations

- Splitting/Fusion can either be enabled/disabled; SLOTHY is not able to determine dynamically if applying the heuristic would be beneficial
  - → State explosion
- Splitting/Fusion is not guaranteed to succeed, e.g., in case of the transformation requiring more registers
- The developer of a callback is responsible for the safety and security of the transformation
  - → Not covered by SLOTHY's selfcheck

## Splitting & Fusion: Limitations

- Splitting/Fusion can either be enabled/disabled; SLOTHY is not able to determine dynamically if applying the heuristic would be beneficial
  - $\rightarrow$ State explosion
- Splitting/Fusion is not guaranteed to succeed, e.g., in case of the transformation requiring more registers
- The developer of a callback is responsible for the safety and security of the transformation
  - $\rightarrow$ Not covered by SLOTHY's selfcheck
- **Current Limitation**: Splitting/Fusion can only enabled/disabled globally.

## Assignment 4: Adding Instructions to Architecture & Microarchitecture Models

### Your Task
Add support for the AArch64 `eon` (Exclusive OR NOT) instruction to SLOTHY

- Operation: `Xd = Xa XOR NOT(Xb)`
- Example: `eon x2, x2, x10`

## Assignment 4: Adding Instructions to Architecture & Microarchitecture Models

### Your Task
Add support for the AArch64 `eon` (Exclusive OR NOT) instruction to SLOTHY

- Operation: `Xd = Xa XOR NOT(Xb)`
- Example: `eon x2, x2, x10`

### Steps

1. **Architecture Model:** Teach SLOTHY to parse `eon` instructions
   `slothy/targets/aarch64/aarch64_neon.py`
2. **Microarchitecture Model:** Add performance characteristics for Cortex-A55
   `slothy/targets/aarch64/cortex_a55.py`

## Assignment 4: Adding Instructions to Architecture & Microarchitecture Models

### Your Task
Add support for the AArch64 `eon` (Exclusive OR NOT) instruction to SLOTHY

- Operation: `Xd = Xa XOR NOT(Xb)`
- Example: `eon x2, x2, x10`

### Steps

1. **Architecture Model:** Teach SLOTHY to parse `eon` instructions
   `slothy/targets/aarch64/aarch64_neon.py`
2. **Microarchitecture Model:** Add performance characteristics for Cortex-A55
   `slothy/targets/aarch64/cortex_a55.py`

### Getting Started

- See files in `04instruction/`
- Read the `README.md` for detailed instructions and hints

## Assignment 4: Setup Instructions

### Clone SLOTHY

- This assignment requires modifying SLOTHY source code
  $\implies$ We need to work on a local clone of SLOTHY
- The assignment code checks that you are not accidentally using SLOTHY from pip

### Setup

1. `cd 04instruction/`
2. `git clone https://github.com/slothy-optimizer/slothy.git`
3. `python3 -m venv venv`
4. `source venv/bin/activate`
5. `pip3 install -r slothy/requirements.txt`
6. Test: `python3 optimize.py` (should fail with parsing error)

## Assignment 4: `test_eon.s`

### Optimization Region (`test_eon.s`)

```
1  slothy_start:
2  ldp x2, x3, [x0]
3  ldp x4, x5, [x0, #16]
4  ldp x19, x20, [x1]
5  ldp x21, x22, [x1, #16]
6
7  eon x2, x2, x19
8  eon x3, x3, x20
9  eon x4, x4, x21
10 eon x5, x5, x22
11
12 stp x2, x3, [x0]
13 stp x4, x5, [x0, #16]
14 slothy_end:
```

### Note

- Simple test with 4 EON operations
- You don't need to modify the assembly

`optimize.py`

```python
1  # Initialize SLOTHY
2  slothy = Slothy(AArch64_Neon, Target_CortexA55)
3
4  # Load the test assembly file
5  slothy.load_source_from_file("test_eon.s")
6
7  # Optimize between markers
8  slothy.optimize(start="slothy_start", end="slothy_end")
9
10 # Write optimized output
11 slothy.write_source_to_file("test_eon_optimized.s")
12
```

**Note:** You do not need to modify `optimize.py` - it will work once EON is added

## Assignment 4: Summary & Bonus Challenge

Your Two Steps

1. **Architecture Model:** Add EON parsing to `aarch64_neon.py`
2. **Microarchitecture Model:** Add EON performance data to `cortex_a55.py`

## Assignment 4: Summary & Bonus Challenge

### Your Two Steps

1. **Architecture Model:** Add EON parsing to `aarch64_neon.py`
2. **Microarchitecture Model:** Add EON performance data to `cortex_a55.py`

### Bonus Challenge - Prize Available!

- First person at CHES to open a PR adding EON to upstream SLOTHY wins a bottle of Taiwanese Whisky!
- Requirements:
    - Must pass CI
    - Must support Cortex-A72 and Neoverse N1 models too
    - Add EON instruction to `tests/naive/aarch64/instructions.s` for CI testing

## Assignment 4: Summary & Bonus Challenge

### Your Two Steps

1. **Architecture Model:** Add EON parsing to `aarch64_neon.py`
2. **Microarchitecture Model:** Add EON performance data to `cortex_a55.py`

### Bonus Challenge - Prize Available!

- First person at CHES to open a PR adding EON to upstream SLOTHY wins a bottle of Taiwanese Whisky!
- Requirements:
    - Must pass CI
    - Must support Cortex-A72 and Neoverse N1 models too
    - Add EON instruction to `tests/naive/aarch64/instructions.s` for CI testing

### Other bonus assignments

- Also support w-form EON: `eon w0, w1, w2`
- Also support Barrel-shifted EON: `eon x0, x1, x2, lsl #8`

## Assignment 5: Instruction Fusion

### Input Code: `fusion.s`

```
1 // Computing column parity in Keccak theta step
2 start:
3     eor v10.16b, v0.16b, v5.16b    // C[0] = A[0,0] ^ A[1,0]
4     eor v10.16b, v10.16b, v11.16b  // C[0] = C[0] ^ A[2,0]
5     eor v10.16b, v10.16b, v16.16b  // C[0] = C[0] ^ A[3,0]
6     eor v10.16b, v10.16b, v21.16b  // C[0] = C[0] ^ A[4,0]
7 end:
```

### Dependencies to Notice

- Chain dependency: Each `eor` depends on previous result
- Target: Use `eor3` instruction from AArch64 SHA3 extension

**Goal**: Write `optimize.py` using SLOTHY's fusion capabilities

## `slothy.fusion_region(start, end, ssa=True)`

Apply fusion callbacks to straightline code region.

- **Input:** `start`, `end` - region labels
- **Output:** Transforms code in-place using fusion callbacks
- **Note:** `ssa`: Output is in static single-assignment (SSA) form (symbolic registers) - usually useful as fusion is used as a pre-processing step before optimization
- **Usage:** `slothy.fusion_region("start", "end", ssa=False)`

# Fusion API Functions

## `slothy.fusion_region(start, end, ssa=True)`

Apply fusion callbacks to straightline code region.

- **Input:** `start, end` - region labels
- **Output:** Transforms code in-place using fusion callbacks
- **Note:** `ssa`: Output is in static single-assignment (SSA) form (symbolic registers) - usually useful as fusion is used as a pre-processing step before optimization
- **Usage:** `slothy.fusion_region("start", "end", ssa=False)`

## `slothy.fusion_loop(loop_lbl, ssa=True)`

Apply fusion callbacks to loop body.

- **Input:** `loop_lbl` - loop label
- **Output:** Transforms loop body using fusion callbacks

[40 minutes hands-on exercise (Assignment 4 + 5)]
See the `README.md` for hints





**Tutorial Assignments**
github.com/dop-amin/ches2025-slothy-tutorial

**Tutorial Slides**
kannwischer.eu/talks/20250914_slothy.pdf

Step 1: Add EON to `aarch64_neon.py`
Add this class (e.g., after the existing `eor` class):

```
class eon(AArch64Instruction):
    pattern = "eon <Xd>, <Xa>, <Xb>"
    inputs = ["Xa", "Xb"]
    outputs = ["Xd"]
```

# Assignment 4: Performance Data from ARM SWOG

## Cortex-A55 Software Optimization Guide

| Instruction group | AArch64 instructions | Exec latency | Execution throughput | Dual-issue | Notes |
|---|---|---|---|---|---|
| ALU, basic, include flag setting | ADD{S}, ADC{S}, AND{S}, BIC{S}, EON, EOR, ORN, ORR, SUB{S}, SBC{S} | 1 | 2 | 11 | – |
| ALU, extend and/or shift | ADD{S}, AND{S}, BIC{S}, EON, EOR, ORN, ORR, SUB{S} | 2 | 2 | 11 | – |

- Basic EON: Latency = 1, Throughput = 2 (2 ALUs, each inverse throughput = 1)
- Shifted EON: Latency = 2, Throughput = 2
- Dual-issue = 11 (both issue slots = can dual issue)

## Assignment 4: Solution - Microarchitecture Model

Add EON to `cortex_a55.py` in 3 dictionaries

```
1  # 1. Execution Unit
2  execution_units = {
3      eor: ExecutionUnit.SCALAR(),
4      eon: ExecutionUnit.SCALAR(), # <- Add this - ALU0 or ALU1
5      # ExecutionUnit.SCALAR() short for [ExecutionUnit.ALU0, ExecutionUnit.ALU1]
6      ...
7  # 2. Throughput (per execution unit!)
8  inverse_throughput = {
9      eor: 1,
10     eon: 1,  # <- Add this
11     ...
12 # 3. Latency
13 latencies = {
14     eor: 1,
15     eon: 1,  # <- Add this
16     ...
```

98

## Architecture Model - Using AArch64Logical Base Class

```
1  # Option 1: Direct inheritance (shown earlier)
2  class eon(AArch64Instruction):
3      pattern = "eon <Xd>, <Xa>, <Xb>"
4      inputs = ["Xa", "Xb"]
5      outputs = ["Xd"]
6
7  # Option 2: Using class hierarchy (more elegant)
8  class eon(AArch64Logical):
9      pattern = "eon <Xd>, <Xa>, <Xb>"
10     inputs = ["Xa", "Xb"]
11     outputs = ["Xd"]
12
```

**Benefit:** Only need one entry in Microarchitecture Model for all instructions

## Assignment 5: Solution

### EOR3 Fusion Solution

```python
1   from slothy import Slothy
2   import slothy.targets.aarch64.aarch64_neon as AArch64_Neon
3   import slothy.targets.aarch64.cortex_a55 as Target_CortexA55
4
5   def main():
6       slothy = Slothy(AArch64_Neon, Target_CortexA55)
7
8       # Load the source assembly file
9       slothy.load_source_from_file("fusion.s")
10
11      # Configure optimization parameters
12      slothy.config.outputs = ["v10"]
13
14      # Perform fusion
15      slothy.fusion_region(start="start", end="end", ssa=False)
16
17      # Write optimized code to output file
18      slothy.write_source_to_file("fusion_opt_a55.s")
```

## Common problem 1: Missing instructions from architecture model

### Error Message

```
ERROR:root:Failed to parse instruction ldr q0, [x1, #16]!
ERROR:root:A list of attempted parsers and their exceptions follows.
...
  File "slothy/targets/aarch64/aarch64_neon.py", line 792, in parser
    raise Instruction.ParsingException(
slothy.targets.aarch64.aarch64_neon.Instruction.ParsingException:
Couldn't parse ldr q0, [x1, #16]!
You may need to add support for a new instruction (variant)?
```

### Solution

- Check if instruction variant is supported in architecture model
- Add new instruction parser to `aarch64_neon.py` if needed

SLOTHY models are built lazily - it is expected that not all instructions are supported.
AArch64 coverage is fairly good by now.

## Common problem 2: Missing instructions from microarchitecture model

### Error Message

```
INFO:slothy.slothy_start.slothy:Attempt optimization with max 0 stalls...
Traceback (most recent call last):
...
  File "slothy/targets/aarch64/cortex_a55.py", line 639, in get_inverse_throughput
    return lookup_multidict(inverse_throughput, src)
  File "slothy/targets/aarch64/aarch64_neon.py", line 4815, in lookup_multidict
    raise UnknownInstruction(f"Couldn't find {instclass} for {inst}")
slothy.targets.common.UnknownInstruction:
Couldn't find <class 'slothy.targets.aarch64.aarch64_neon.vmls'> for mls v11.4S, v12.4S, v13.4S
```

### Solution

- Instruction is parsed correctly but missing performance data in microarchitecture model
- Add execution unit assignment and latency data to target model (e.g., `cortex_a55.py`)
- Check software optimization guide for the target microarchitecture

SLOTHY models are built incrementally - not all instructions have complete microarchitecture data yet.

# Common problem 3: Trying to optimize non-straight-line code

## Code

```
1   start:
2     ldr  q0, [x1], #16
3     add  v1.4s, v0.4s, v0.4s
4     str  q1, [x2], #16
5
6     subs x3, x3, #1
7     b.ne start  // <-- Problem!
8
9     mov  v2.16b, v1.16b
10  end:
```

Trying to use: `slothy.optimize(start="start", end="end")`

## Error Message

```
ERROR:root:Failed to parse instruction b.ne start
...
slothy.targets.aarch64.aarch64_neon
  .Instruction.ParsingException:
Couldn't parse b.ne start
You may need to add support for a new
instruction (variant)?
```

## Solution

- SLOTHY cannot optimize across branches
- For loops: use `optimize_loop("start")`
- For other branches: optimize each basic block separately
- Place markers to avoid branch instructions

# Common problem 4: Useless instruction warnings

## Loop Code

```
1   loop_start:
2     ldr  q0, [x0], #16
3     ldr  q1, [x1], #16
4
5     add  v2.4s, v0.4s, v1.4s
6
7     str  q2, [x2], #16
8
9     subs x4, x4, #1
10    b.ne loop_start
```

## Error Message

```
ERROR:slothy.loop_start.slothy.dataflow:
The result registers ['x0'] of instruction
0:[ldr q0, [x0], #16] are neither used nor
declared as global outputs.
ERROR:slothy...dataflow:This is often a configuration
  error. Did you miss an output declaration?
  ...
slothy.core.dataflow.SlothyUselessInstructionException:
Useless instruction detected
```

## Solution

- Use: `config.inputs_are_outputs = True` for loops
- Use: `config.outputs = [...] for other outputs`
- Warnings: `config.allow_useless_instructions` should not be used in those cases
- Note: Can also be caused by incorrect input/output declarations in the arch. model

# Common problem 5: SLOTHY using callee-saved registers

## Original Function

```
1   .global vec_add
2   vec_add:
3   start:
4       ldr q0, [x0]
5       ldr q1, [x1]
6       add v2.4s, v0.4s, v1.4s
7       str q2, [x0]
8   end:
9       ret
```

## SLOTHY Output (BROKEN!)

```
1   .global vec_add
2   vec_add:
3   start:
4       ldr q7, [x0]
5       ldr q3, [x1]
6       add v8.4s, v7.4s, v3.4s
7       str q8, [x0]
8   end:
9       ret
```

### Problem

- SLOTHY renamed v2 → v8
- v8-v15 are callee-saved
- Function corrupts v8

### Solution

```
1   # Reserve callee-saved vector regs
2   slothy.config.reserved_regs = [
3       "v8", "v9", "v10", "v11",
4       "v12", "v13", "v14", "v15"
5   ]
```

# Thank you for joining the SLOTHY tutorial!

### Get Involved:

- Found a bug or have a feature request? Open an issue on GitHub!
- Want to contribute? Pull requests are always welcome!
- We have plenty of ideas for further research
- Questions? Feel free to reach out directly

### Contact:

Matthias Kannwischer matthias@kannwischer.eu

Amin Abdulrahman amin.abdulrahman@mpi-sp.org

**GitHub:** github.com/slothy-optimizer/slothy

Join our Discord!
discord.gg/Khy2bwgm