

Predictable Accelerator Design with Time-Sensitive Affine Types

Background:

- Specialized hardware to deal with diminishing Moore's law returns
- Reconfigurable hardware: FPGAs are hard to program
- Designed in Verilog, VHDL, Bluespec, Chisel
RTL: Register transfer level
Describes how data + instructions travel in the circuit at a painful detail
- Vendors offer "HLS" solutions

C to gates

 Verilog / VHDL

Comment: Verilog/VHDL is irreplaceable for performance/power/area tuning,

LUTs in FPGA

Problem with HLS tools:

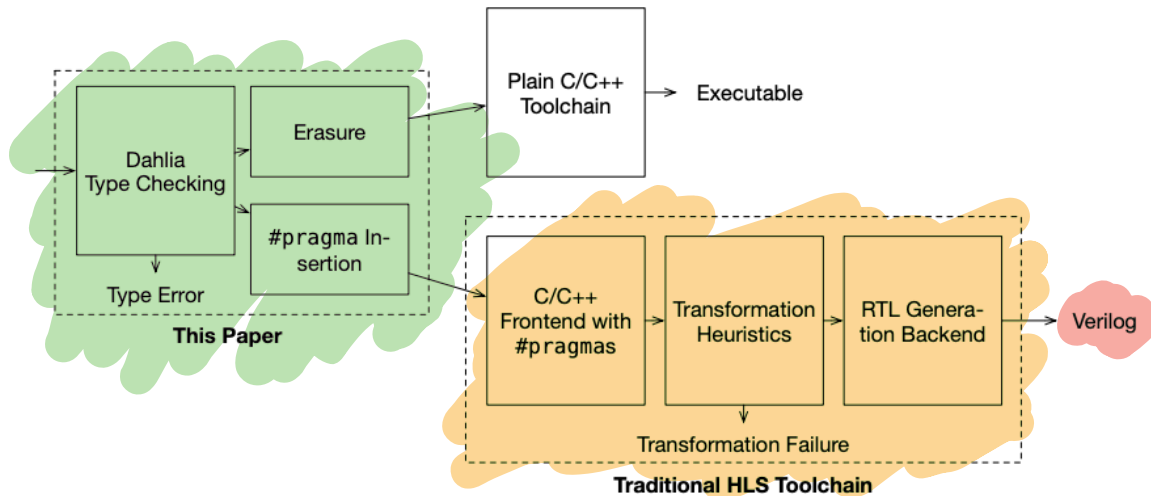
- Support only a subset of C
- Heuristics based approach to generating hardware
- Minor, Smoother changes result in large swings

Unpredictable hardware generation

Solution:

- Restrict to HLS programs with clear hardware implementation
- Explicit annotations for "costly" implementations.

Overview :



Predictability by enforcing *Subunits* to do one thing at a time.
Substructural logics are great at that!
linear, Affine

Previous works with similar idea

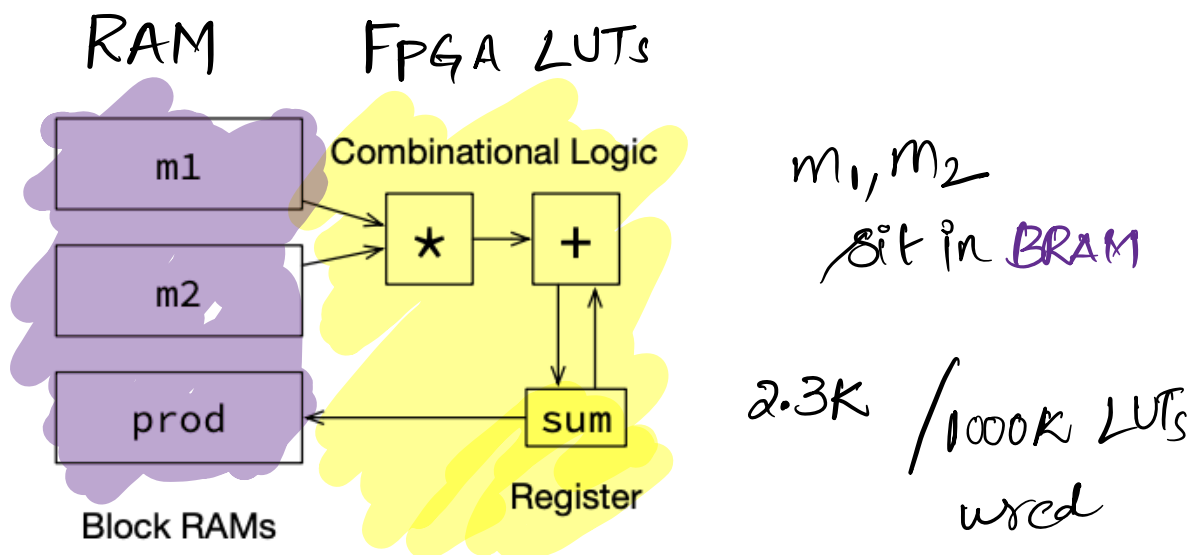
- Rust!
- Region-based mem allocation in cyclone
PLDI 2002
- A type & effect system for deterministic
parallel Java: COPSLA
2000

Unpredictability of HLS tools:

```
1 int m1[512][512], m2[512][512], prod[512][512];
2 int sum;
3 for (int i = 0; i < 512; i++) {
4     for (int j = 0; j < 512; j++) {
5         sum = 0;
6         for (int k = 0; k < 512; k++) {
7             sum += m1[i][k] * m2[k][j];
8         }
9         prod[i][j] = sum; } }
```

Figure 2. Dense matrix multiplication in HLS-friendly C.

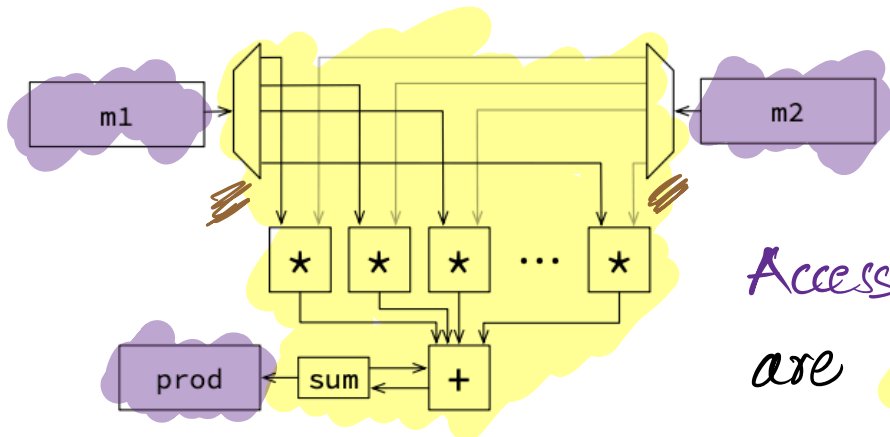
Verilog would have clocks, ports and resets



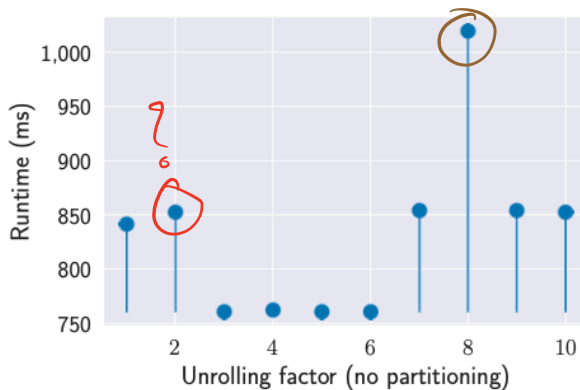
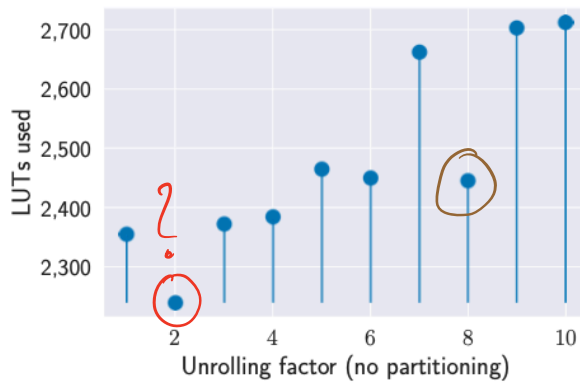
(a) The original code.

#Pragma HLS UNROLL FACTOR=8

```
for (int k = 0; k < 512; k++) {  
    sum += m1[i][k] * m2[k][j];  
}
```



(b) With unrolling.



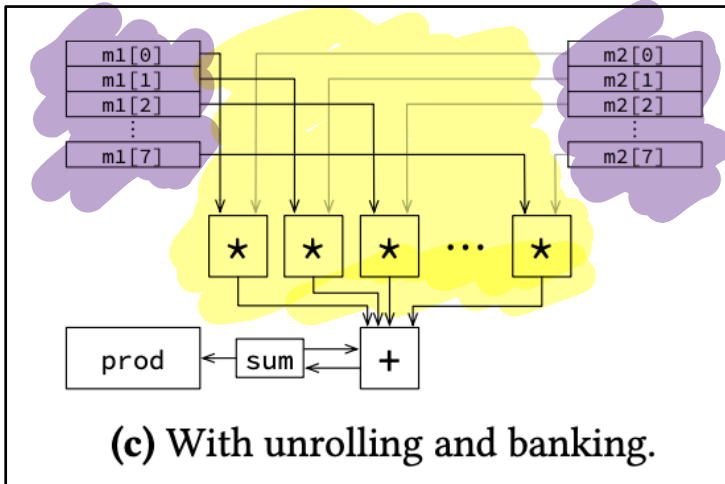
Muxing is the bottleneck.

HLS ARRAY-PARTITION VAR = m_1 FACTOR = 8

* HLS ARRAY-PARTITION VAR = m_2 FACTOR = 8

Vars split across

BRAMS



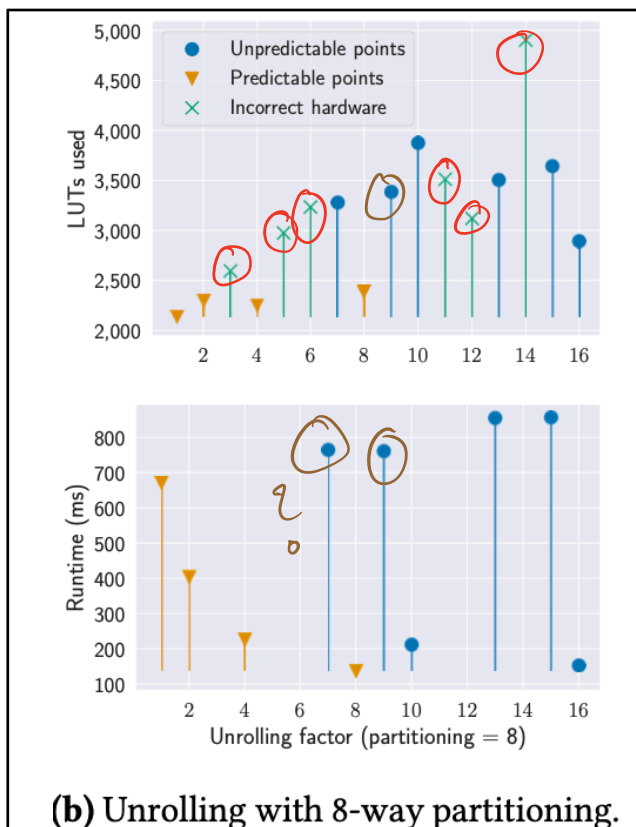
$m[0] \rightarrow \text{bank}[0]$

$m[1] \rightarrow \text{bank}[1]$

\vdots

$m[7] \rightarrow \text{bank}[7]$

$m[8] \rightarrow \text{bank}[0]$

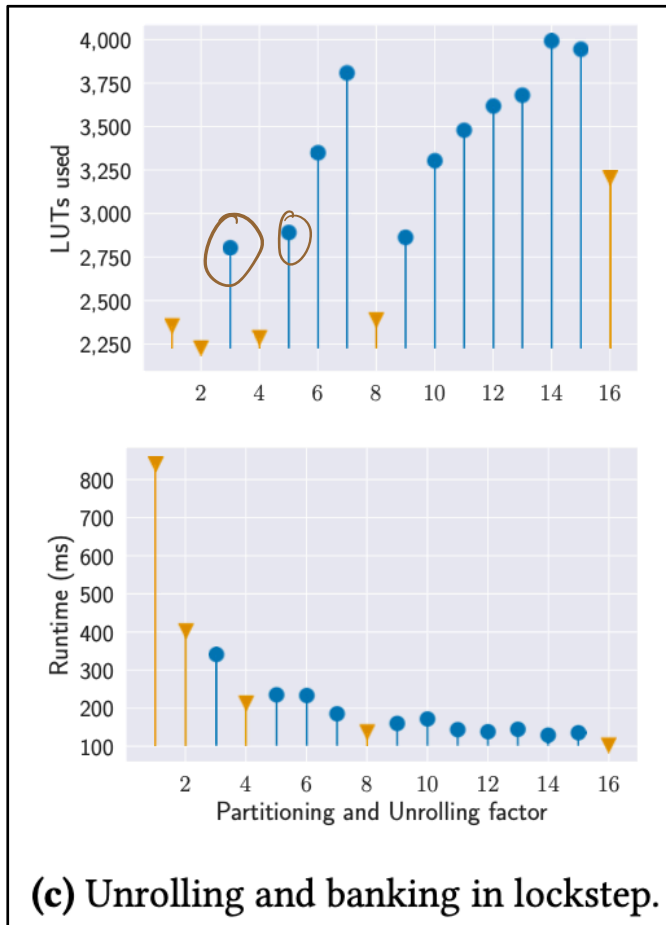


Bank Parallelism

doesn't scale

monotonically!

Both unrolling and banking in lockstep



Good performance
when $\text{bank} \% \text{unroll} = 0$

2, 4, 8, 16

Otherwise requires
complicated

Muxing to
access data

Ex: bank factor = 8

unroll factor = 9

Every combinatorial circuit needs to
access every bank

Dahlia language:

Safety property imposed

Simultaneous accesses to a given bank
may not exceed the number of ports.

```
let A: float[10];
```

BRAM in FPGA

// memories defined by type
and size.

✓

```
let x = A[5];
```

✗

```
let B = A;
```

No aliasing or
duplication of mems.

✓

```
let x = A[0]; // OK  
A[1] := 1;    // Error:
```

✗

The read x consumed A

✓

```
let x = A[0];
```

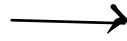
✓

```
let y = A[0];
```

Same address simultaneous
reads are OK!

Reads are non-affine when on
same index of same mem


```
let x = A[0];
let y = A[0];
```



```
let tmp = A[0];
let x = tmp;
let y = tmp;
```

" ; "

unordered composition: statements execute in parallel

" --- " ordered composition: sequential

✓

```
let x = A[0]
---
A[1] := 1
```

Read and write are sequential.

Restoration of resource access.

```
let A: float[10]; let B: float[10];
{
  let x = A[0] + 1
  ---
  B[1] := A[1] + x // OK
};
let y = B[0]; // Error: B already consumed.
```

sequential

Parallel

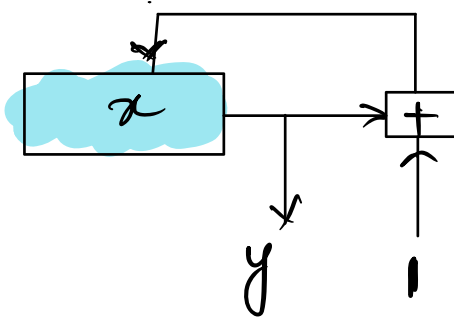
Abstract notion of time, not tied to underlying design's clock cycles.

```
let x = a + b; let y = c * d
```

Local vars : wires & registers

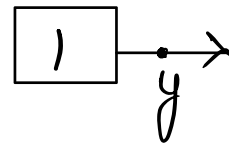
✓ `let x = 0; x := x + 1; let y = x;`

no affine typing here



if x is needed
across clock cycles

\gg



if x is unused
elsewhere

`let x = A[0] + 1 --- B[0] := A[1] + x`

x is a register in this case.

Memory banking support

`let A: Float [n bank m];`

`// n % m = 0`

Explicit affine tracking for each bank



```
let A: float[10 bank 2];  
A{0}[0] := 1;  
A{1}[0] := 2; // OK: Accessing
```

$\text{float}\{2\}[10]$
Alt Syntax

```
let M: float{2} [4] {2} [4];  
// [4 bank 2] [4 bank 2]
```

$2 \times 2 = 4$ banks.

0,0		0,1	
0	1	2	3
4	5	6	7
8	9	A	B
C	D	E	F
1,0		1,1	

loops, unrolling:

```
for (let i = 0..10) unroll 2 { f(i) }
```

\equiv

```
for (let i = 0..5) { f(2*i + 0); f(2*i + 1) }
```

for is parallelizable, but no cross-iteration dependencies.

X

```
let A: float[10];  
for (let i = 0..10) unroll 2 {  
  A[i] := compute(i) // Error: Insufficient banks.  
}
```

for (let $i = 0 \dots 8$) unroll 4 { $A[i]$ }

\uparrow
idx { 0..4 }

\nwarrow
check for 4 access
resources

Interaction b/w composition & unrolling

```
let A: float[10 bank 2];  
for (let i = 0..10) unroll 2 {  
  let x = A[i]  
  ---  
  f(x, A[0]) }
```

```
for (let i = 0..5) {  
  { let  $x_0$  = A[2*i] --- f( $x_0$ , A[0]) };  
  { let  $x_1$  = A[2*i + 1] --- f( $x_1$ , A[0]) } }
```

Parallelization within each logical timestep

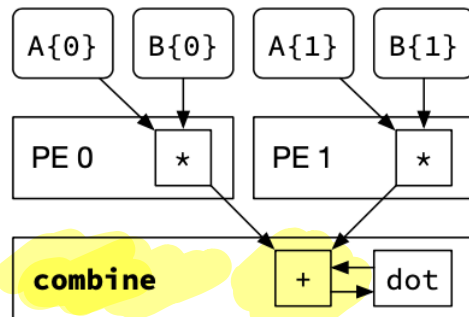
```
for (let i = 0..5) {  
  { let x0 = A[2*i]; let x1 = A[2*i + 1] }  
  ---  
  { f(x0, A[0]); f(x1, A[0]) } }
```

Combine construct:

```
for (let i = 0..10) unroll 2 { dot += A[i] * B[i] }
```

$+=$ introduces silent cross iteration dependency
 $\text{dot} += A[2i] * B[2i]; \text{dot} += A[2i+1] * B[2i+1];$

```
for (let i = 0..10)  
unroll 2 {  
  let v = A[i] * B[i];  
} combine {  
  dot += v;  
}
```



Sequential code to be performed after each unrolled iteration

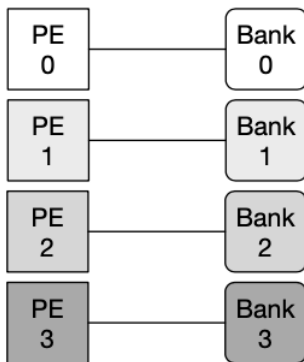
\rightarrow is a combine register
 \rightarrow tuple of all values
 $+=, -=, /=, *=$

Memory views for flexible iteration;

Arbitrary computations in indices result in bad weird hardware. Ex: $A[i * i]$

Dahlia restricts index expr to simple expressions
 $[i + 8]$

Hence: logical arrangements.

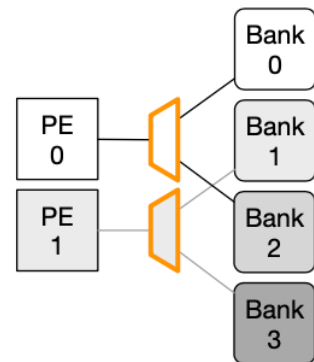


direct
access

```
let A: float [8 bank 4]
for (let i = 0..8) unroll 4
  A[i];
```

Shrink:

```
let A: float [8 bank 4];
view sh = shrink A[by 2]; // sh: float [8 bank 2]
for (let i = 0..8) unroll 2
  sh[i]; // OK: sh has 2 banks. Compiled to: A[i].
```



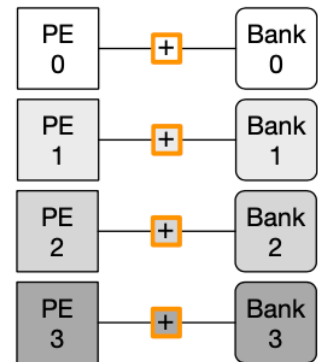
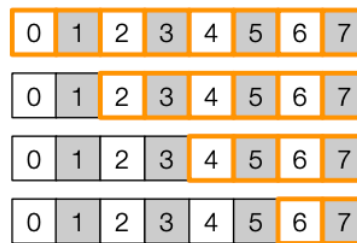
Suffix:

```
view v = suffix M[by k * e];
```

$V_{\{b\}}[i] \equiv M_{\{b\}}[i + e]$

$k =$ bank factor of M .

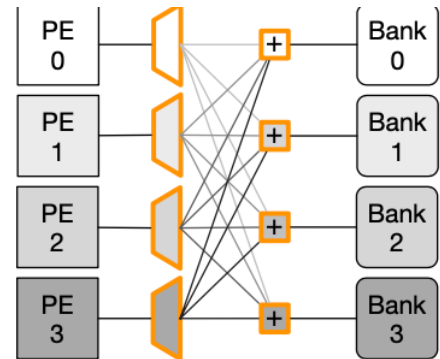
```
let A: float[8 bank 2];
for (let i = 0..4) {
  view s = suffix A[by 2*i];
  s[1]; // reads A[2*i + 1]
}
```



Shift:

```
view v = shift M[by e];
```

```
let A: float[12 bank 4];
for (let i = 0..3) {
  view r = shift A[by i*i]; // r: float[12 bank 4]
  for (let j = 0..4) unroll 4
    let x = r[j]; // accesses A[i*i + j]
}
```



Split:

For efficient h/w gen in nested loops.

```
let A, B: float[12 bank 4];
view shA, shB = shrink A[by 2], B[by 2];
for (let i = 0..6) unroll 2 {
  view vA, vB = suffix shA[by 2*i], shB[by 2*i];
  for (let j = 0..2) unroll 2 {
    let v = vA[j] + vB[j];
  } combine {
    sum += v; }}}
```

Dahlia cannot reason about separation of v_A, v_B

```
float A[12], B[12], sum = 0.0;
for (int i = 0; i < 6; i++)
  for (int j = 0; j < 2; j++)
    sum += A[2*i + j] * B[2*i + j];
```

$\text{view } v_A = \text{split } A \text{ [by 2]}$

Formalism:

$x \in \text{variables} \quad a \in \text{memories} \quad n \in \text{numbers}$
 $b ::= \text{true} \mid \text{false} \quad v ::= n \mid b$
 $e ::= v \mid \text{bop } e_1 \ e_2 \mid x \mid a[e]$
 $c ::= e \mid \text{let } x = e \mid c_1 \text{ — } c_2 \mid c_1 ; c_2 \mid \text{if } x \ c_1 \ c_2 \mid$
 $\quad \text{while } x \ c \mid x := e \mid a[e_1] := e_2 \mid \text{skip}$
 $\tau ::= \text{bit}\langle n \rangle \mid \text{float} \mid \text{bool} \mid \text{mem } \tau[n_1]$

for compiles to
while

Large Step Semantics:

$$\frac{a \notin \rho_1 \quad \sigma_1, \rho_1, e \Downarrow \sigma_2, \rho_2, n \quad \sigma_2(a)(n) = v}{\sigma_1, \rho_1, a[e] \Downarrow \sigma_2, \rho_2 \cup \{a\}, v}$$

mem access

\checkmark : Vars to mem map
 \int : Set of mems. accessed.

$$\frac{\sigma_1, \rho_1, c_1 \Downarrow \sigma_2, \rho_2 \quad \sigma_2, \rho_2, c_2 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, c_1 ; c_2 \Downarrow \sigma_3, \rho_3}$$

Parallel
composition

$$\frac{\sigma_1, \rho_1, c_1 \Downarrow \sigma_2, \rho_2 \quad \sigma_2, \rho_1, c_2 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, c_1 \text{ — } c_2 \Downarrow \sigma_3, \rho_2 \cup \rho_3}$$

Sequential
composition

mems accessed in c_1 are freed back (ρ_1)

Type system:

Γ : variable context (for checking index, local types)

Δ : Affine context for memories.

$$\frac{\Gamma, \Delta_1 \vdash e_1 : \text{bit}(n) \dashv \Delta_2 \quad \Delta_2 = \Delta_3 \cup \{a \mapsto \text{mem } \tau[n_1]\}}{\Gamma, \Delta_1 \vdash a[e] : \tau \dashv \Delta_3}$$

mem
access

why the Δ_2 update for index?

$$\frac{\Gamma_1, \Delta_1 \vdash c_1 \dashv \Gamma_2, \Delta_2 \quad \Gamma_2, \Delta_2 \vdash c_2 \dashv \Gamma_3, \Delta_3}{\Gamma_1, \Delta_1 \vdash c_1 ; c_2 \dashv \Gamma_3, \Delta_3}$$

Parallel composition

$$\frac{\Gamma_1, \Delta_1 \vdash c_1 \dashv \Gamma_2, \Delta_2 \quad \Gamma_2, \Delta_1 \vdash c_2 \dashv \Gamma_3, \Delta_3}{\Gamma_1, \Delta_1 \vdash c_1 \text{ --- } c_2 \dashv \Gamma_3, \Delta_2 \cap \Delta_3}$$

seq composition

Lemma 1 (Progress). If $\Gamma, \Delta \vdash c \dashv \Gamma_2, \Delta_2$ and $\Gamma, \Delta \sim \sigma, \rho$, then $\sigma, \rho, c \rightarrow \sigma', \rho', c'$ or $c = \text{skip}$.

Lemma 2 (Preservation). If $\Gamma, \Delta \vdash c \dashv \Gamma_2, \Delta_2$ and $\Gamma, \Delta \sim \sigma, \rho$, and $\sigma, \rho, c \rightarrow \sigma', \rho', c'$, then $\Gamma', \Delta' \vdash c' \dashv \Gamma'_2, \Delta'_2$ and $\Gamma', \Delta' \sim \sigma', \rho'$.

Theorem. If $\emptyset, \Delta^* \vdash c \dashv \Gamma_2, \Delta_2$ and $\emptyset, \emptyset, c \xrightarrow{*} \sigma, \rho, c'$ and $\sigma, \rho, c' \not\rightarrow$, then $c' = \text{skip}$.

Results:

Experiments on mach HLS suite: 16/19

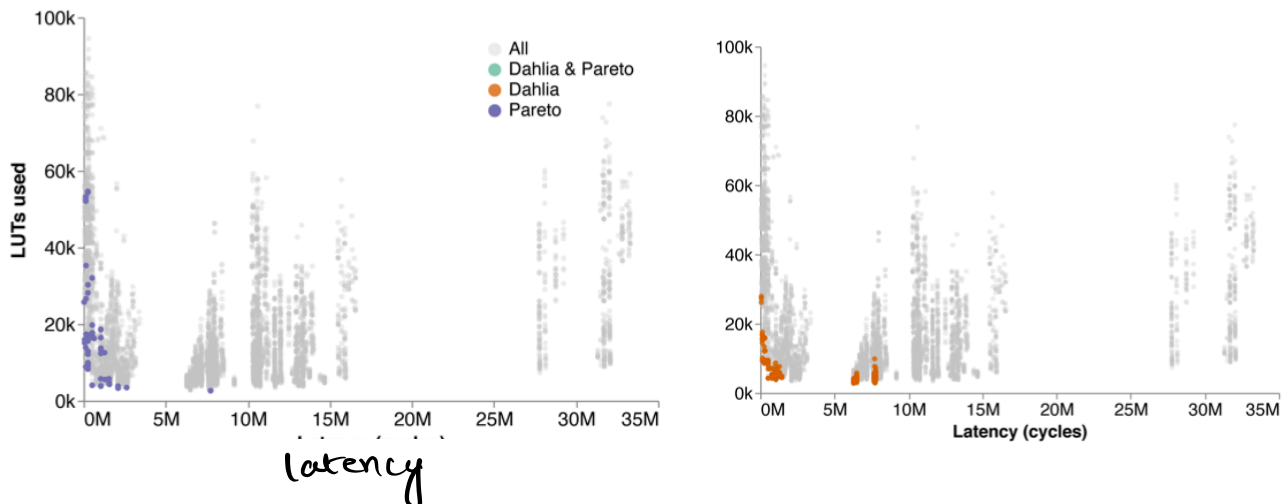
successful

```
for (r=0; r<row_size-2; r++)  
  for (c=0; c<col_size-2; c++)  
    for (k1=0; k1<3; k1++)  
      for (k2=0; k2<3; k2++)  
        mul = filter[k1*3 + k2] *  
              orig[(r+k1)*col_size + c+k2];
```

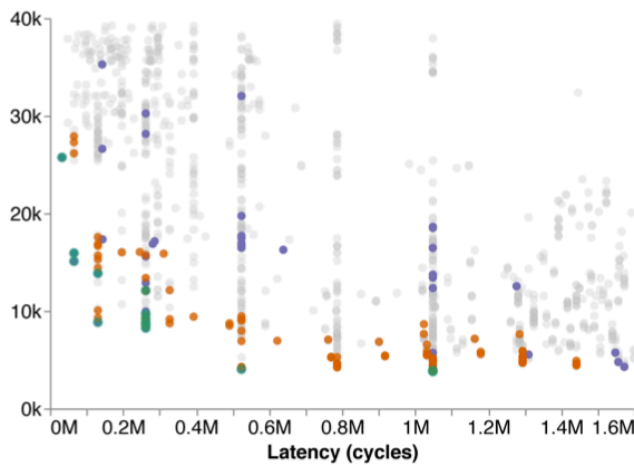
C code

```
for (let row = 0..126) {  
  for (let col = 0..62) {  
    view window = shift orig[by row][by col];  
    for (let k1 = 0..3) unroll 3 {  
      for (let k2 = 0..3) unroll 3 {  
        let mul = filter[k1][k2] * window[k1][k2];
```

Dahlia code



On a macro level the programs accepted are near "optimal".



However not the exact points that traditional tools & dahlia accept.

Things that went right :

1. Deterministic design generation
2. Abstract view of time*
3. Substructural logics in H/W, an interesting idea

Things that went wrong: the

1. C-like paradigm might not be right mental model for H/W in general.
2. The hardware level optimisations manifest as weird syntax
3. Unlike S/W langt, user needs to inspect the generated h/w and that ability isn't clear in general with HLS

4. verilog has generate & for that
behave similarly !!