

# BITSAD v2: Compiler Optimization and Analysis for Bitstream Computing

HiPEAC 2020 (published in TACO)

---

Kyle Daruwalla, Heng Zhuo, Rohit Shukla, and Mikko Lipasti

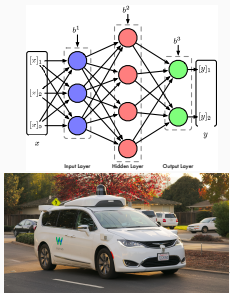
20 January 2020

University of Wisconsin - Madison, Dept. of Elec. and Comp. Eng.

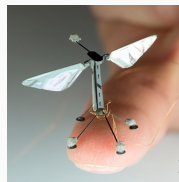


# Motivation

CV/ML algorithms enable powerful new applications



Fabrication techniques enabling pico-aerial vehicles (PAVs)

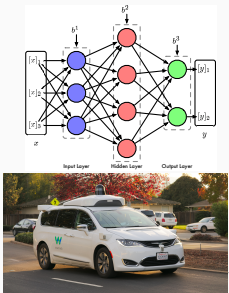


<sup>1</sup>Dllu 2017

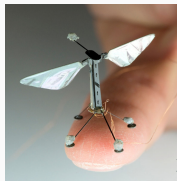
<sup>2</sup>Ma 2015

# Motivation

CV/ML algorithms enable powerful new applications

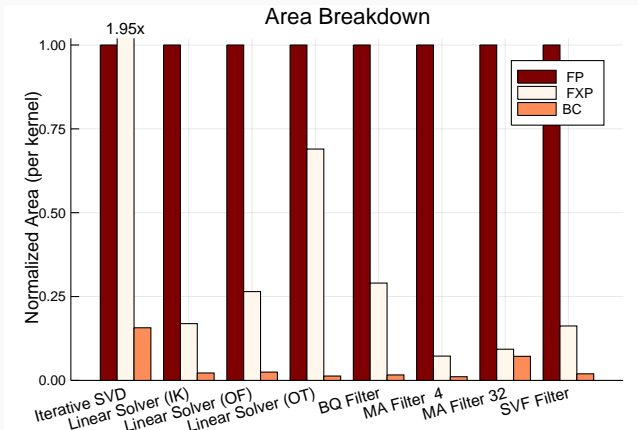


Fabrication techniques enabling pico-aerial vehicles (PAVs)



How do we enable advanced CV/ML algorithms on ultra-low power platforms?

# Motivation



The resource consumption of bitstream computing (BC) implementations is much lower than floating point (FP) and fixed point (FXP) designs.

1. Background: Bitstreams and example applications
2. BITSAD: What is it good for? <sup>3</sup>
3. Population coding: Parallelization for stochastic computing
4. Optimizations
  - 4.1 for stochastic bitstreams
  - 4.2 for deterministic bitstreams

---

<sup>3</sup>Seinfeld: S5E14, Feb. 1994.

## Background: Bitstream Computing

---

# What is Bitstream Computing?

## Stochastic Bitstreams:



$$\begin{aligned} \mathbb{E}[S_1] &= 0.5 \frac{0, 1, 0, 1, 1, 1, 0, 0}{1, 1, 0, 1, 1, 0, 1, 1} \\ \mathbb{E}[S_2] &= 0.75 \frac{1, 1, 0, 1, 1, 0, 1, 1}{1, 1, 0, 1, 1, 0, 1, 1} \end{aligned} \quad \text{AND} \quad \begin{aligned} \mathbb{E}[S_1 S_2] &= \frac{0, 1, 0, 1, 1, 0, 0, 0}{1, 1, 0, 1, 1, 0, 1, 1} \\ &= \mathbb{E}[S_1] \mathbb{E}[S_2] \\ &= 0.375 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[S_1] &= 0.5 \frac{0, 1, 0, 1, 1, 1, 0, 0}{1, 0, 0, 0, 1, 0, 0, 0} \\ \mathbb{E}[S_2] &= 0.25 \frac{1, 0, 0, 0, 1, 0, 0, 0}{1, 0, 0, 0, 1, 0, 0, 0} \end{aligned} \quad \text{OR} \quad \begin{aligned} \mathbb{E}[S_1 + S_2] &= \frac{1, 1, 0, 1, 1, 1, 0, 0}{1, 0, 0, 0, 1, 0, 0, 0} \\ &= \mathbb{E}[S_1] + \mathbb{E}[S_2] \\ &= 0.75 \end{aligned}$$

# What is Bitstream Computing?

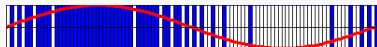
## Stochastic Bitstreams:



$$\begin{aligned} \mathbb{E}[S_1] &= 0.5 \frac{0, 1, 0, 1, 1, 1, 0, 0}{1, 1, 0, 1, 1, 0, 1, 1} \\ \mathbb{E}[S_2] &= 0.75 \frac{1, 1, 0, 1, 1, 0, 1, 1}{1, 1, 0, 1, 1, 0, 1, 1} \end{aligned} \quad \text{AND} \quad \begin{aligned} &0, 1, 0, 1, 1, 0, 0, 0 \\ &= \mathbb{E}[S_1] \mathbb{E}[S_2] \\ &= 0.375 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[S_1] &= 0.5 \frac{0, 1, 0, 1, 1, 1, 0, 0}{1, 1, 0, 0, 0, 1, 0, 0} \\ \mathbb{E}[S_2] &= 0.25 \frac{1, 0, 0, 0, 1, 0, 0, 0}{1, 1, 0, 1, 1, 1, 0, 0} \end{aligned} \quad \text{AND} \quad \begin{aligned} &1, 1, 0, 1, 1, 1, 0, 0 \\ &= \mathbb{E}[S_1 + S_2] \\ &= \mathbb{E}[S_1] + \mathbb{E}[S_2] \\ &= 0.75 \end{aligned}$$

## Deterministic Bitstreams:



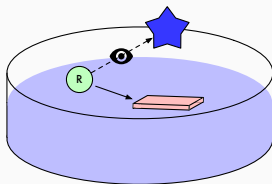
Density of "1"  $\Rightarrow$  Higher  
amplitude

- Sequence is *deterministic*
- Oversampled audio data
- Leads to efficient filters



## Example Application #1: Navigation

Consider a robot (denoted “R”) that needs to navigate an unknown environment

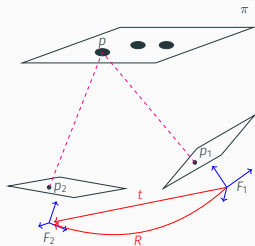
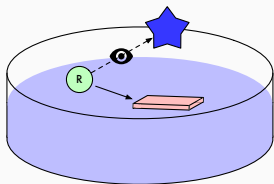


---

<sup>4</sup>Malis and Vargas 2007.

## Example Application #1: Navigation

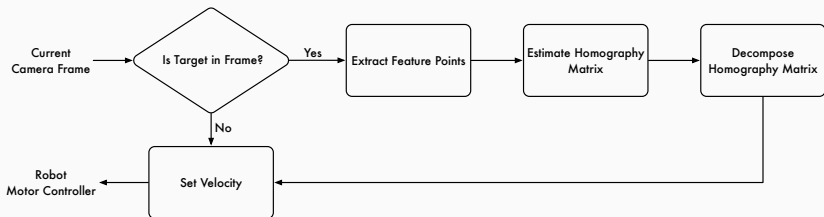
Consider a robot (denoted “R”) that needs to navigate an unknown environment



Navigation by visual cues – homography estimation and decomposition<sup>4</sup>

<sup>4</sup>Malis and Vargas 2007.

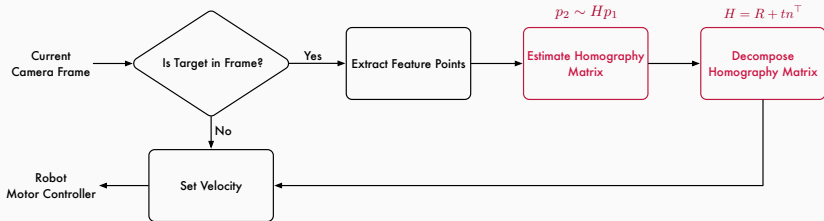
# Example Application #1: Navigation Pipeline



---

<sup>5</sup>Dubrofsky 2009.

# Example Application #1: Navigation Pipeline



Requires linear solver and singular value decomposition<sup>5</sup>

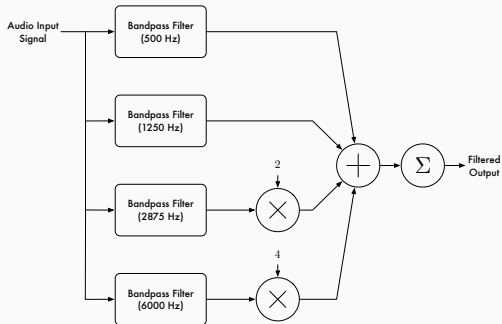
---

<sup>5</sup>Dubrofsky 2009.

## Example Application #2: Channel Shaping

Audio information is sent across several frequency channels

Need to shape input signal to focus on channels of interest

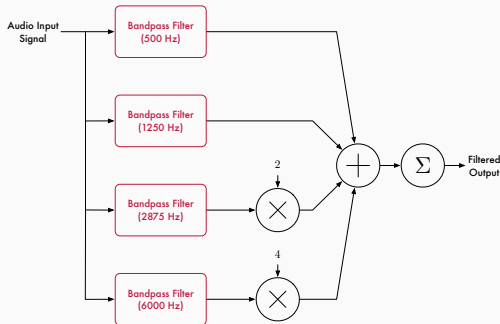


## Example Application #2: Channel Shaping

Audio information is sent across several frequency channels

Need to shape input signal to focus on channels of interest

Oversampled bitstreams  $\Rightarrow$  more efficient filters



BITSAD

---

# What is BITSAD?

A domain-specific language for bitstream computing



**sbt**

Allows users to simulate bit-level, cycle-accurate designs



# What is BITSAD?

A domain-specific language for bitstream computing



**sbt**

Allows users to simulate bit-level, cycle-accurate designs  
*and* automatically generate Verilog implementations!

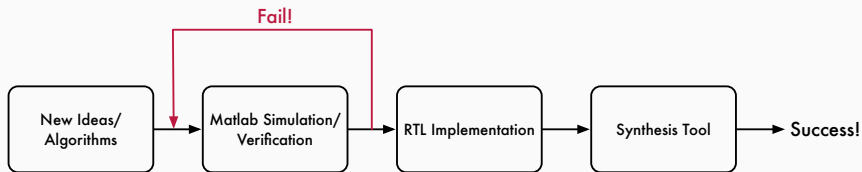
# Typical Design Flow

Traditional design flow turnover is slow:



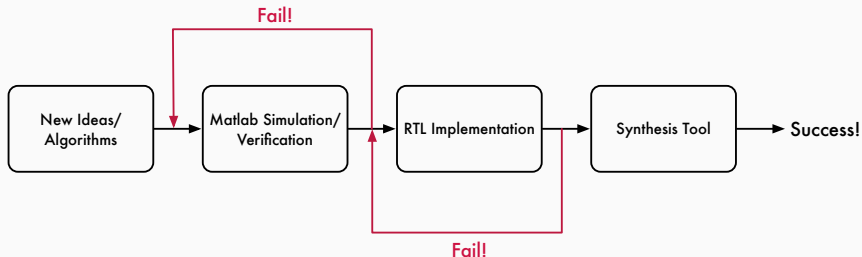
# Typical Design Flow

Traditional design flow turnover is slow:



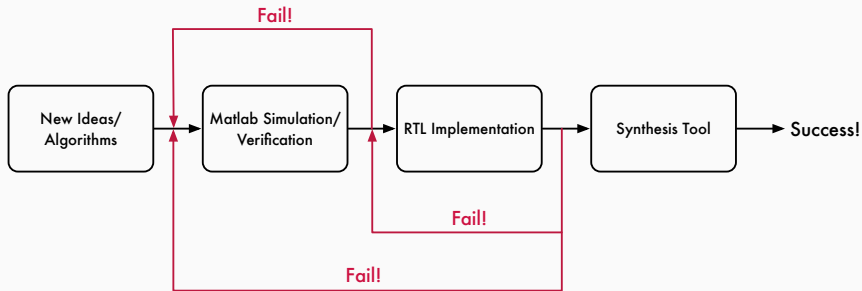
# Typical Design Flow

Traditional design flow turnover is slow:



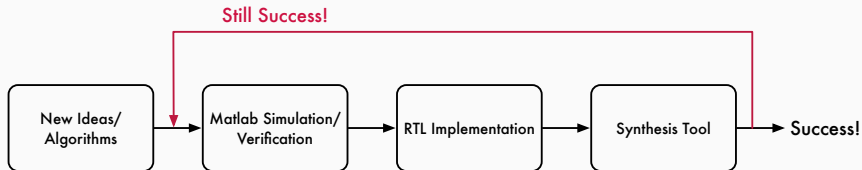
# Typical Design Flow

Traditional design flow turnover is slow:



# Typical Design Flow

Traditional design flow turnover is slow:



# Typical Design Flow

Where the problem is coming from:



# Typical Design Flow

What can we do about it:





---

## Algorithm 1 Iterative SVD

---

**Require:** Input matrix  $A \in \mathbb{R}^{m \times n}$  and initial guess  $v_0 \in \mathbb{R}^n$

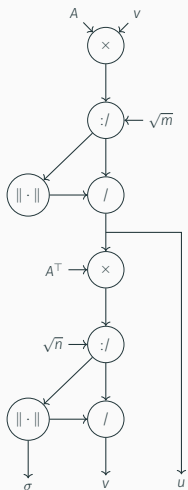
- 1: **for**  $k = 1, 2, \dots$  (until convergence) **do**
  - 2:    $w_k = Av_{k-1}$
  - 3:    $\alpha_k = \|w_k\|_2 = \sqrt{w_k^T w_k}$
  - 4:    $u_k = w_k / \alpha_k$
  - 5:    $z_k = A^T u_k$
  - 6:    $\sigma_k = \|z_k\|_2 = \sqrt{z_k^T z_k}$
  - 7:    $v_k = z_k / \sigma_k$
  - 8: **end for**
  - 9: **return** First left/right singular vectors,  $u_k$  &  $v_k$ , and first singular value,  $\sigma_k$
- 

---

```
1 def loop(A: Matrix[SBitstream],
2         v: Matrix[SBitstream]):
3     (Matrix[SBitstream],
4      Matrix[SBitstream], SBitstream) = {
5
6     // Update right singular vector
7     var w = A * v
8     var wScaled = w ./ math.sqrt(params.m)
9     var u = wScaled / Matrix.norm(wScaled)
10
11    // Update left singular vector
12    var z = A.T * u
13    var zScaled = z ./ math.sqrt(params.n)
14    var sigma = Matrix.norm(zScaled)
15    var _v = zScaled / sigma
16
17    (u, _v, sigma)
18 }
```

---

# BITSAD by Example

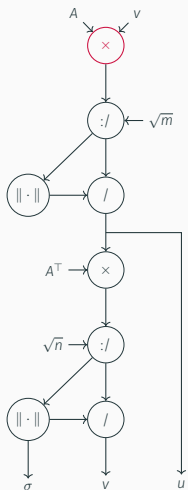


---

```
1 def loop(A: Matrix[SBitstream],
2         v: Matrix[SBitstream]):
3     (Matrix[SBitstream],
4      Matrix[SBitstream], SBitstream) = {
5
6     // Update right singular vector
7     var w = A * v
8     var wScaled = w ./ math.sqrt(params.m)
9     var u = wScaled / Matrix.norm(wScaled)
10
11    // Update left singular vector
12    var z = A.T * u
13    var zScaled = z ./ math.sqrt(params.n)
14    var sigma = Matrix.norm(zScaled)
15    var _v = zScaled / sigma
16
17    (u, _v, sigma)
18 }
```

---

# BITSAD by Example

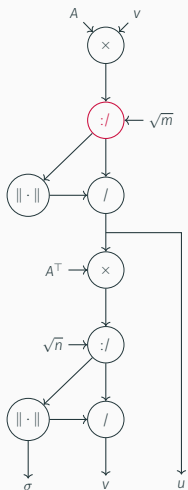


---

```
1 def loop(A: Matrix[SBitstream],
2         v: Matrix[SBitstream]):
3     (Matrix[SBitstream],
4      Matrix[SBitstream], SBitstream) = {
5
6     // Update right singular vector
7     var w = A * v
8     var wScaled = w ./ math.sqrt(params.m)
9     var u = wScaled / Matrix.norm(wScaled)
10
11    // Update left singular vector
12    var z = A.T * u
13    var zScaled = z ./ math.sqrt(params.n)
14    var sigma = Matrix.norm(zScaled)
15    var _v = zScaled / sigma
16
17    (u, _v, sigma)
18 }
```

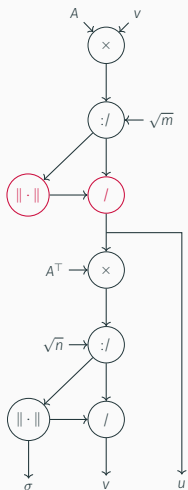
---

# BITSAD by Example



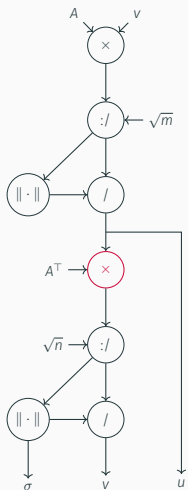
```
1 def loop(A: Matrix[SBitstream],  
2         v: Matrix[SBitstream]):  
3     (Matrix[SBitstream],  
4      Matrix[SBitstream], SBitstream) = {  
5  
6     // Update right singular vector  
7     var w = A * v  
8     var wScaled = w :/ math.sqrt(params.m)  
9     var u = wScaled / Matrix.norm(wScaled)  
10  
11    // Update left singular vector  
12    var z = A.T * u  
13    var zScaled = z :/ math.sqrt(params.n)  
14    var sigma = Matrix.norm(zScaled)  
15    var _v = zScaled / sigma  
16  
17    (u, _v, sigma)  
18 }
```

# BITSAD by Example



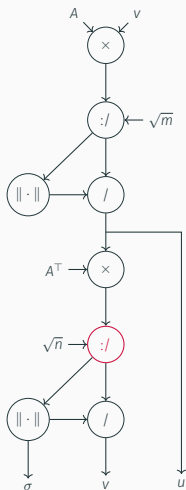
```
1 def loop(A: Matrix[SBitstream],  
2         v: Matrix[SBitstream]):  
3     (Matrix[SBitstream],  
4      Matrix[SBitstream], SBitstream) = {  
5  
6     // Update right singular vector  
7     var w = A * v  
8     var wScaled = w :/ math.sqrt(params.m)  
9     var u = wScaled / Matrix.norm(wScaled)  
10  
11    // Update left singular vector  
12    var z = A.T * u  
13    var zScaled = z :/ math.sqrt(params.n)  
14    var sigma = Matrix.norm(zScaled)  
15    var _v = zScaled / sigma  
16  
17    (u, _v, sigma)  
18 }
```

# BITSAD by Example



```
1 def loop(A: Matrix[SBitstream],  
2         v: Matrix[SBitstream]):  
3     (Matrix[SBitstream],  
4      Matrix[SBitstream], SBitstream) = {  
5  
6     // Update right singular vector  
7     var w = A * v  
8     var wScaled = w :/ math.sqrt(params.m)  
9     var u = wScaled / Matrix.norm(wScaled)  
10  
11    // Update left singular vector  
12    var z = A.T * u  
13    var zScaled = z :/ math.sqrt(params.n)  
14    var sigma = Matrix.norm(zScaled)  
15    var _v = zScaled / sigma  
16  
17    (u, _v, sigma)  
18 }
```

# BITSAD by Example

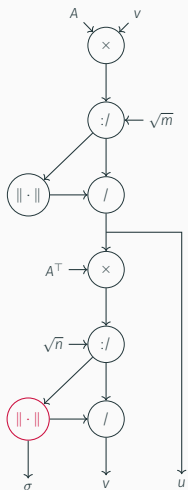


---

```
1 def loop(A: Matrix[SBitstream],
2         v: Matrix[SBitstream]):
3     (Matrix[SBitstream],
4      Matrix[SBitstream], SBitstream) = {
5
6     // Update right singular vector
7     var w = A * v
8     var wScaled = w ./ math.sqrt(params.m)
9     var u = wScaled / Matrix.norm(wScaled)
10
11    // Update left singular vector
12    var z = A.T * u
13    var zScaled = z ./ math.sqrt(params.n)
14    var sigma = Matrix.norm(zScaled)
15    var _v = zScaled / sigma
16
17    (u, _v, sigma)
18 }
```

---

# BITSAD by Example



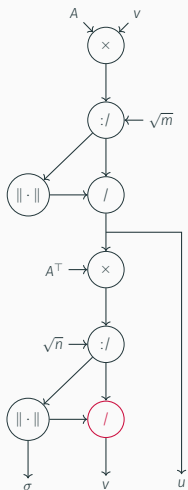
---

```
1 def loop(A: Matrix[SBitstream],
2         v: Matrix[SBitstream]):
3     (Matrix[SBitstream],
4      Matrix[SBitstream], SBitstream) = {
5
6     // Update right singular vector
7     var w = A * v
8     var wScaled = w :/ math.sqrt(params.m)
9     var u = wScaled / Matrix.norm(wScaled)
10
11    // Update left singular vector
12    var z = A.T * u
13    var zScaled = z :/ math.sqrt(params.n)
14    var sigma = Matrix.norm(zScaled)
15    var _v = zScaled / sigma
16
17    (u, _v, sigma)
18 }
```

---



# BITSAD by Example

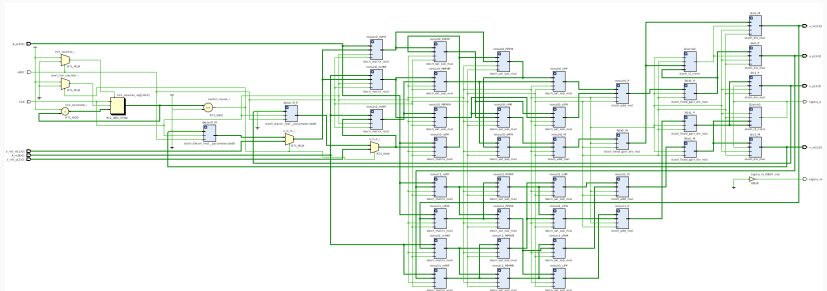


---

```
1 def loop(A: Matrix[SBitstream],
2         v: Matrix[SBitstream]):
3     (Matrix[SBitstream],
4      Matrix[SBitstream], SBitstream) = {
5
6     // Update right singular vector
7     var w = A * v
8     var wScaled = w ./ math.sqrt(params.m)
9     var u = wScaled / Matrix.norm(wScaled)
10
11    // Update left singular vector
12    var z = A.T * u
13    var zScaled = z ./ math.sqrt(params.n)
14    var sigma = Matrix.norm(zScaled)
15    var _v = zScaled / sigma
16
17    (u, _v, sigma)
18 }
```

---

# BITSAD by Example



- Allows designers to write algorithms at a high abstraction level
- Simulates designs with bit-level, cycle-accurate results
- Generates synthesizable hardware with no code changes

# BITSAD Summary

- Allows designers to write algorithms at a high abstraction level
- Simulates designs with bit-level, cycle-accurate results
- Generates synthesizable hardware with no code changes

Is this enough to realize “general-purpose” bitstream computing?

# Population Coding

---

# Latency Issue for Stochastic Bitstreams

Recall the definition of a stochastic bitstream:

$$\frac{1}{T} \sum_{t=1}^T X_t \approx \mathbb{E}X_t = p \quad \text{as } T \rightarrow \infty$$

# Latency Issue for Stochastic Bitstreams

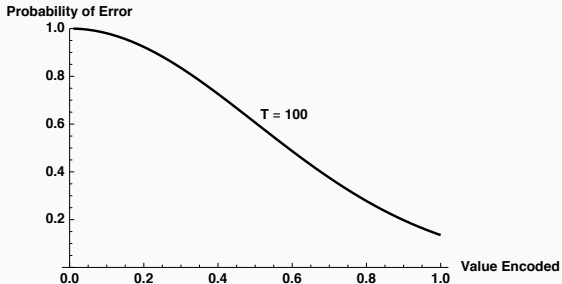
Recall the definition of a stochastic bitstream:

$$\frac{1}{T} \sum_{t=1}^T X_t \approx \mathbb{E}X_t = p \quad \text{as } T \rightarrow \infty$$

How large should  $T$  be for 10% application error? 5% error?

# Latency Issue for Stochastic Bitstreams

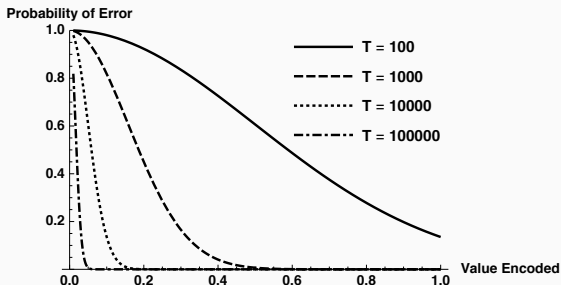
Quantify error using Hoeffding's inequality:





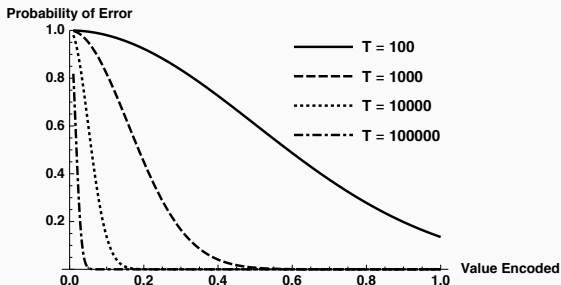
# Latency Issue for Stochastic Bitstreams

Quantify error using Hoeffding's inequality:



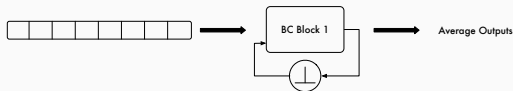
# Latency Issue for Stochastic Bitstreams

Quantify error using Hoeffding's inequality:



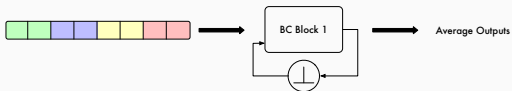
Error depends on # of time steps, relative accuracy, *and* magnitude of number!

# Parallelizing Stochastic Bitstreams



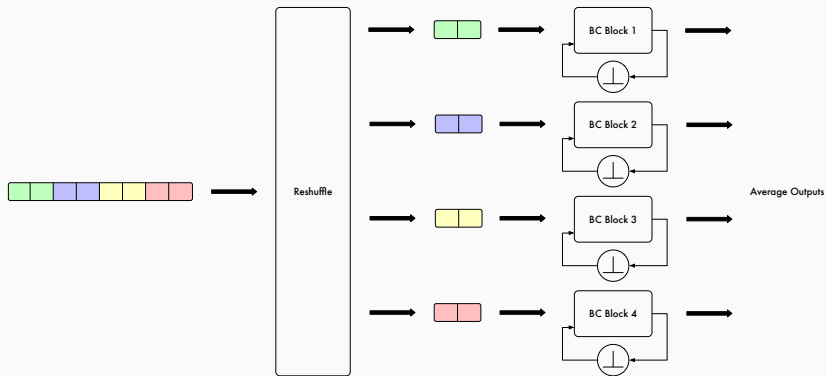
$$p \approx \frac{1}{T} \sum_{t=1}^T X_t$$

# Parallelizing Stochastic Bitstreams



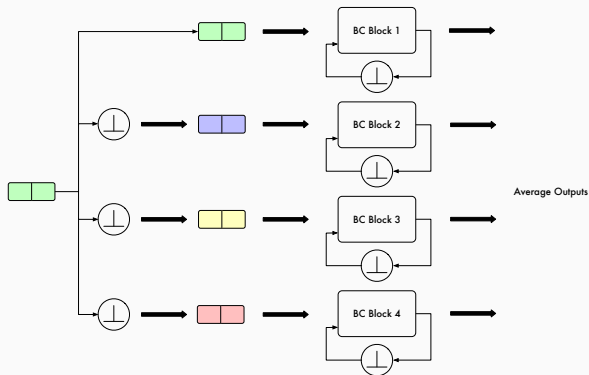
$$p \approx \frac{1}{T} \sum_{t=1}^T X_t$$

# Parallelizing Stochastic Bitstreams



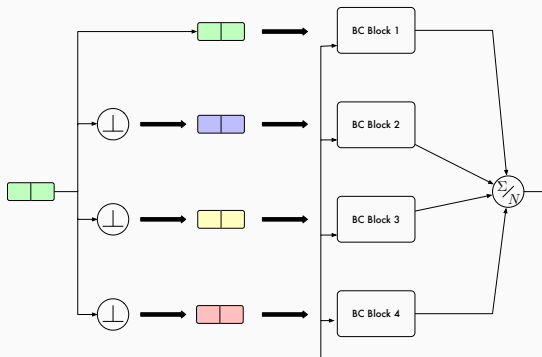
$$p \approx \frac{1}{T} \sum_{t=1}^T X_t$$

# Parallelizing Stochastic Bitstreams



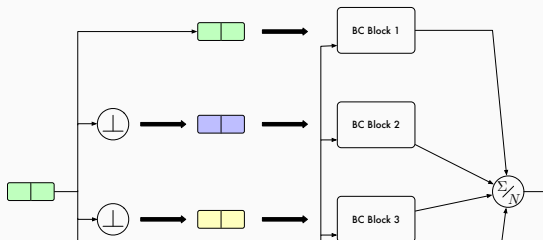
$$p \approx \frac{1}{T/N} \sum_{t=1}^{T/N} \frac{1}{N} \sum_{i=1}^N X_{t,i}$$

# Parallelizing Stochastic Bitstreams



$$p \approx \frac{1}{T/N} \sum_{t=1}^{T/N} \frac{1}{N} \sum_{i=1}^N X_{t,i}$$

# Parallelizing Stochastic Bitstreams



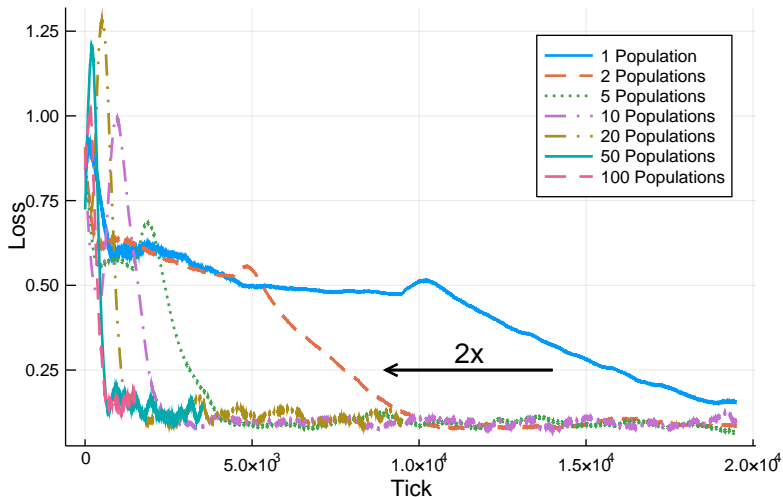
**Theorem:** For all interesting circuits, population coding guarantees that output can be fed back without decorrelation

$$p \approx \frac{1}{T/N} \sum_{t=1}^{T/N} \frac{1}{N} \sum_{i=1}^N X_{t,i}$$



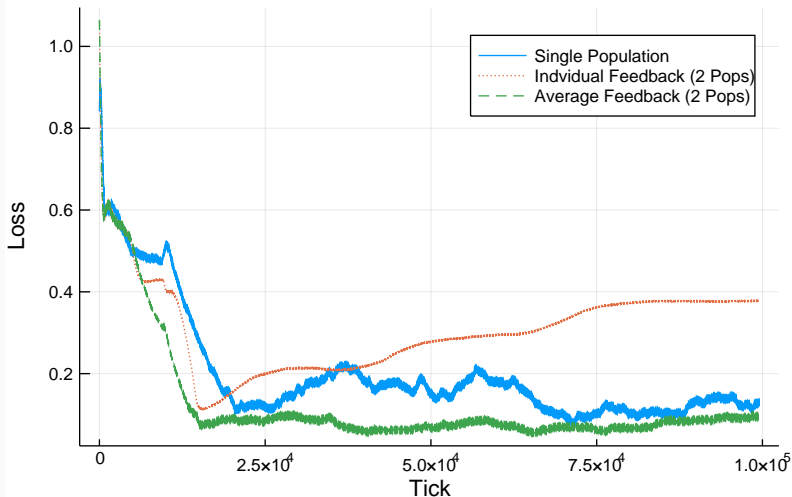
# Population Coding Experiments

## Average Loss for Iterative SVD over 10 Trials



# Population Coding Experiments

## Average Loss for Iterative SVD over 10 Trials



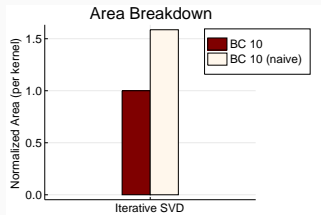
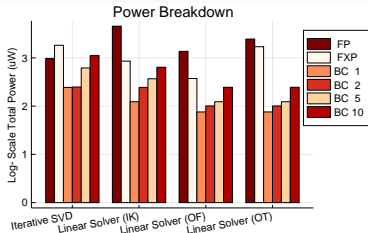
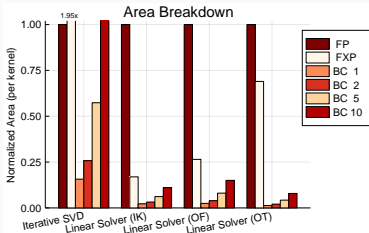
# Population Coding in BITSAD

---

```
1 def loop(A: Matrix[SBitstream], v: Matrix[SBitstream]):
2     (Matrix[SBitstream], Matrix[SBitstream], SBitstream)
3     = populations(2, {
4 // Update right singular vector
5 var w = A * v
6 var wScaled = w ./ math.sqrt(params.m)
7 var u = wScaled / Matrix.norm(wScaled)
8
9 // Update left singular vector
10 var z = A.T * u
11 var zScaled = z ./ math.sqrt(params.n)
12 var sigma = Matrix.norm(zScaled)
13 var _v = zScaled / sigma
14
15 (u, _v, sigma)
16 }
```

---

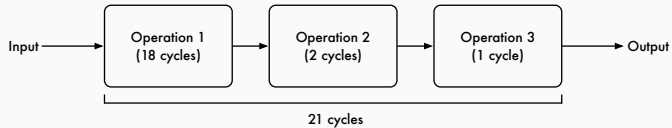
# Population Coding Results



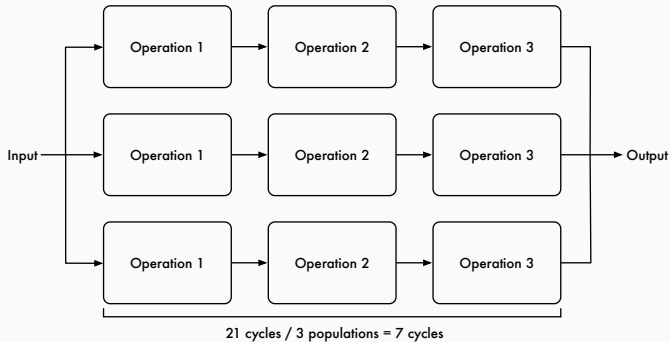
# Optimizations for Stochastic Bitstreams

---

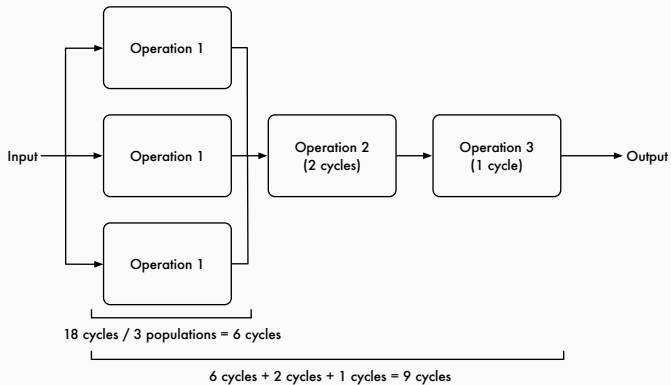
# Population Coding Granularity



# Population Coding Granularity

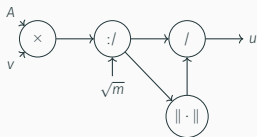


# Population Coding Granularity





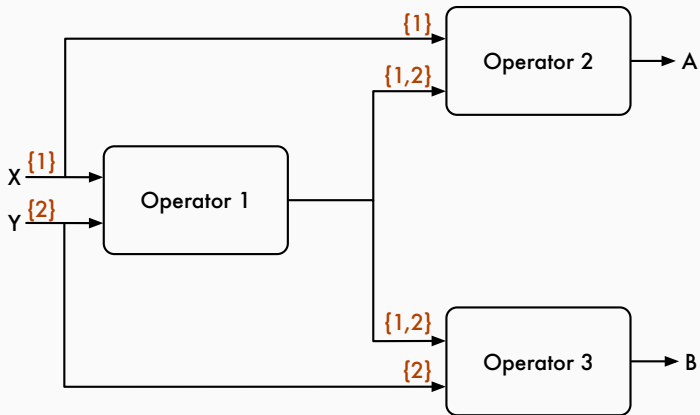
# Preliminary Results of Granular Pop Coding



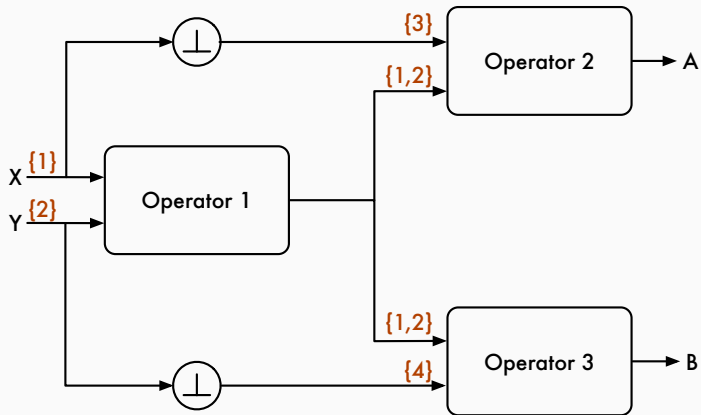
**Table 1:** A Comparison of Designs with Population Coding at Different Granularities

Variant	Average Cycles till Convergence	Area (# LUTs + # FFs)
8 Populations (Full Design)	19781	6807
8 Populations (Multiplier Only)	13613	4484

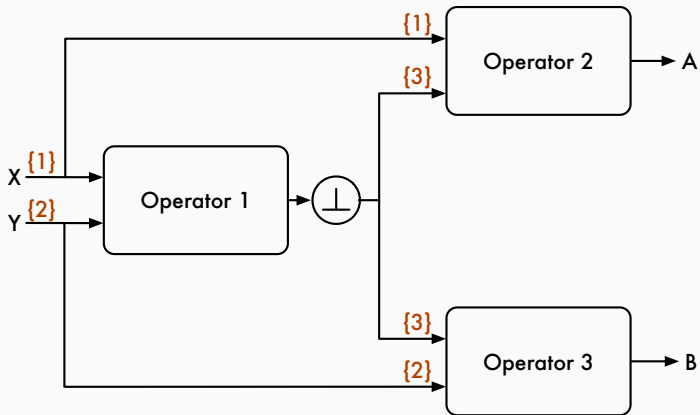
## Optimal Placement of Decorrelators



# Optimal Placement of Decorrelators



# Optimal Placement of Decorrelators



# Optimizations for Deterministic Bitstreams

---

# Deterministic Bitstream Optimization

Consider the following digital state-variable filter:

---

```
1// Get delay buffer values
2val d1_old = delay1.pop
3val d2_old = delay2.pop
4
5// Update SDM outputs (f and q are compile time constants)
6val d2 = sdm2.evaluate(f * d1_old + d2_old)
7val d1 = sdm1.evaluate(f * (x - d2 - q * d1_old) + d1_old)
8
9// Push new values into delay buffers
10delay1.push(d1)
11delay2.push(d2)
```

---

# Deterministic Bitstream Optimization

Consider the following digital state-variable filter:

---

```
1// Get delay buffer values
2val d1_old = delay1.pop
3val d2_old = delay2.pop
4
5// Update SDM outputs (f and q are compile time constants)
6val d2 = sdm2.evaluate(f * d1_old + d2_old)
7val d1 = sdm1.evaluate(f * (x - d2 - q * d1_old) + d1_old)
8
9// Push new values into delay buffers
10delay1.push(d1)
11delay2.push(d2)
```

---

With strength reduction:

---

```
1val d2 = sdm2.evaluate(f * d1_old + d2_old)
2val d1 = sdm1.evaluate(f * x - f * d2 - (f * q) * d1_old + d1_old)
```

---

# Deterministic Bitstream Optimization

Consider the following digital state-variable filter:

---

```
1// Get delay buffer values
2val d1_old = delay1.pop
3val d2_old = delay2.pop
4
5// Update SDM outputs (f and q are compile time constants)
6val d2 = sdm2.evaluate(f * d1_old + d2_old)
7val d1 = sdm1.evaluate(f * (x - d2 - q * d1_old) + d1_old)
8
9// Push new values into delay buffers
10delay1.push(d1)
11delay2.push(d2)
```

---

With strength reduction:

---

```
1val d2 = sdm2.evaluate(f * d1_old + d2_old)
2val d1 = sdm1.evaluate(f * x - f * d2 - (f * q) * d1_old + d1_old)
```

---

With algebraic simplification:

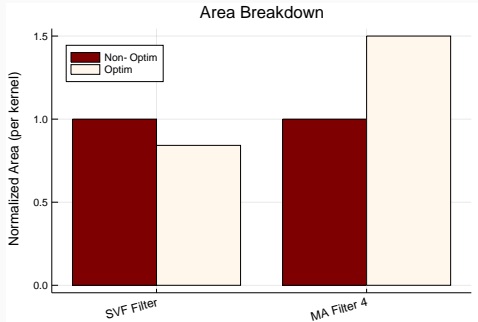
---

```
1val d2 = sdm2.evaluate(f * d1_old + d2_old)
2val d1 = sdm1.evaluate(f * x - f * d2 + (1 - f * q) * d1_old)
```

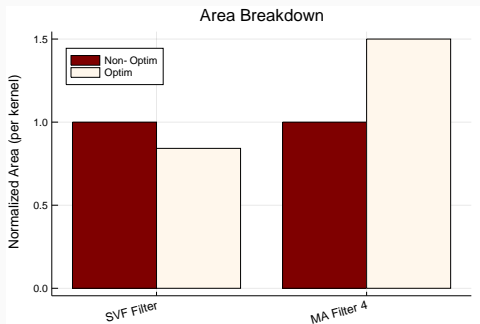
---



# Deterministic Bitstream Optimization



# Deterministic Bitstream Optimization



Synthesize submodules with standard cell library and provide area estimates to compiler as a clue

## Concluding Remarks

---

# Conclusion

Our language, BITSAD, allows users to:

1. design algorithms at a high level
2. simulate bit-level, cycle-accurate results
3. generate hardware automatically

We also introduced:

1. population coding to parallelize stochastic computing circuits w/o sacrificing accuracy
2. optimizations for deciding population coding granularity and decorrelator placement
3. optimizations for deterministic bitstream designs

Questions?

Check out BITSAD on GitHub:

<https://github.com/UW-PHARM/BitSAD>

<https://github.com/UW-PHARM/BitBench>

## References i



Dllu (2017). URL: [https://commons.wikimedia.org/wiki/File:Waymo\\_Chrysler\\_Pacifica\\_in\\_Los\\_Altos,\\_2017.jpg](https://commons.wikimedia.org/wiki/File:Waymo_Chrysler_Pacifica_in_Los_Altos,_2017.jpg).



Dubrofsky, Elan (2009). "Homography Estimation". PhD thesis. Carleton University.



Ma, Kevin Y. (2015). *RoboBee*. URL: <http://www.aboutkevinma.com/index.html#publications> (visited on 04/01/2018).



Malis, Ezio and Manuel Vargas (2007). “Deeper understanding of the homography decomposition for vision-based control”. In: *Sophia* 6303.6303, p. 90. ISSN: 0036-8075. DOI: 10.1126/science.318.5857.1691b. URL: <http://hal.archives-ouvertes.fr/inria-00174036/>.