#### The Anytime Automaton

Joshua San Miguel Natalie Enright Jerger



#### Summary

#### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.

#### Summary

#### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.





#### Summary

#### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.



#### Many applications are inherently noisy and imprecise.



#### Many applications are inherently noisy and imprecise.



6



```
program ()
{
   foos_on_first();
   bars_on_second();
   hello_worlds_on_third();
```





}

UNIVERSITY OF TORONTO



UNIVERSITY OF TORONTO



```
program ()
 {
      approx_foos_on_first();
      approx_bars_on_second();
      hello_worlds_on_third();
 }
                                             hello_worlds_on_third
  foos_on_first
                        bars_on_second
time
                                       tune quality
```

```
program ()
 {
      approx_foos_on_first();
      approx_bars_on_second();
      hello_worlds_on_third();
 }
                                        hello_worlds_on_third
                     bars_on_second
  foos_on_first
time
                                  tune quality
```

UNIVERSITY OF TORONTO



12

of Electrical & Computer Engineering

**NIVERSITY OF TORONTO** 

Difficult to ensure acceptability of final output on-the-fly, since quality control limited to local approximations and not their composition. (Challenge #1: Holistic Quality Control)



Difficult to ensure acceptability of final output on-the-fly, since quality control limited to local approximations and not their composition. (Challenge #1: Holistic Quality Control)



Real-time systems impose strict runtime constraints; loss in output quality more tolerable than not finishing in time.



```
program ()
 {
      approx_foos_on_first();
      approx_bars_on_second();
      hello_worlds_on_third();
 }
                                                  hello_worlds_dn_third
    foos_on_first
                             bars_on_second
time
                                                        strict target runtime
```







Difficult to ensure strict real-time constraints are met (i.e., interrupt the application), since runtime-quality tradeoffs vary dynamically. (Challenge #2: Interruptibility)



In user-interactive environments, users dictate quality requirements on-the-fly.



```
program ()
{
    approx_foos_on_first();
    approx_bars_on_second();
    hello_worlds_on_third();
```





}





Difficult to ensure acceptability for a given user at a given context, since acceptable quality cannot be determined a priori.



#### We propose the **Anytime Automaton**:

- A new computation model for approximate computing.
- Revisits and generalizes concepts from anytime (or iterative) algorithms, originally studied for real-time decision problems.
- A recipe for applying approximate computing techniques such that the final output is available early and improves in quality over time.

quality



#### application
















# Outline

#### **Anytime Automaton**

- The Model
- The Approximations

#### Evaluation

- Methodology
- Experimental Results

#### Conclusion







The Edward S. Rogers Sr. Department of Electrical & Computer Engineering UNIVERSITY OF TORONTO

**.**9,



The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
UNIVERSITY OF TORONTO





The Edward S. Rogers Sr. Department of Electrical & Computer Engineering UNIVERSITY OF TORONTO







































# **Data Diffusion Model**











































































# Anytime Automaton – The Model

1. Ensure precise output is always produced eventually.

## Anytime Automaton – The Model

1. Ensure precise output is always produced eventually.



## Anytime Automaton – The Model

1. Ensure precise output is always produced eventually.


1. Ensure precise output is always produced eventually.



1. Ensure precise output is always produced eventually.



2. Create the effect of improving accuracy over time.

Anytime (or iterative) algorithms have been studied before but are traditionally built into the *coarse-grained* derivation of an application.

Approximate computing techniques have proliferated recently and have been shown to have general *fine-grained* applicability.









2. Create the effect of improving accuracy over time.



The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
UNIVERSITY OF TORONTO

3. Enable interruptibility via pipelining.

3. Enable interruptibility via pipelining.



3. Enable interruptibility via pipelining.





3. Enable interruptibility via pipelining.





# Anytime Automaton – The Model 3. Enable interruptibility via pipelining. computation C time final output not ready!



3. Enable interruptibility via pipelining.



3. Enable interruptibility via pipelining.



3. Enable interruptibility via pipelining.





3. Enable interruptibility via pipelining.



1. General case: apply approximations iteratively.

1. General case: apply approximations iteratively.



quality knob

1. General case: apply approximations **iteratively**.

loop perforation	smaller perforation stride
floating-point precision	more mantissa bits
SRAM bit upsets	higher supply voltage
load value approximation	lower approximation degree
neural acceleration	higher neural network complexity

- 1. General case: apply approximations iteratively.
- Loop perforation

- 1. General case: apply approximations iteratively.
- Loop perforation

 $\bigcirc$ 

perforation stride 20: for i = 0, 20, 40, 60, 80, 100, ...., N-1

General case: apply approximations iteratively.
Loop perforation

$\bigcirc$	perforation stride 20:	for i = 0, 20, 40, 60, 80, 100,, N-1
$\bigcirc$	perforation stride 15:	for i = 0, 15, 30, 45, 60, 75,, N-1
$\bigcirc$	perforation stride 10:	for i = 0, 10, 20, 30, 40, 50,, N-1
	perforation stride 5:	for i = 0, 5, 10, 15, 20, 25,, N-1
	perforation stride 1:	for i = 0, 1, 2, 3, 4, 5, 6, 7,, N-1

# General case: apply approximations iteratively. Loop perforation

perforation stride 20: for i = 0, 20, 40, 60, 80, 100, ...., N-1

# Achieves desired effect of improving quality over time, but can yield redundant work.

perforation stride 5: for i = 0, 5, 10, 15, 20, 25, ...., N-1

**perforation stride 1:** for i = 0, 1, 2, 3, 4, 5, 6, 7, ..., N-1

General case: apply approximations iteratively.
Loop perforation



- 2. Better case: apply **diffusive** approximations.
- > Each approximation builds on the previous one.

- 2. Better case: apply **diffusive** approximations.
- > Each approximation builds on the previous one.



data dependences (each approximate result contributes usefully to precise result)

- 2. Better case: apply **diffusive** approximations.
- > Each approximation builds on the previous one.



data dependences (each approximate result contributes usefully to precise result)

data sampling	more samples	
		/
nteger/fixed-point precision	more bits	

- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)



- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)



- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)



- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)

To improve quality, no need to reiterate from beginning; therefore, *diffusive*. (e.g., just add more samples to current result)

- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)



- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)



- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)


- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)



- 2. Better case: apply **diffusive** approximations.
- Input sampling (e.g., generating a distribution)

Minimal redundant work since each element processed exactly once.

- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)



- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)



- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)



- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)



- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)





- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)







- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)





- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)



- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)





- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)





- 2. Better case: apply **diffusive** approximations.
- Output sampling (e.g., generating an image)





- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)

#### [ x y z ] • [ 10.1101 01.0010 11.0110 ]

- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)

#### [ X Y Z ] • [ 10.1101 01.0010 11.0110 ]



- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)





- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)

#### [ X Y Z ] • [ 10.1101 01.0010 11.0110 ]

X * 10.1101	Y * 01.0010	Z * 11.0110
		$\rightarrow$

- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)



- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)



- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)



- 2. Better case: apply **diffusive** approximations.
- Integer/fixed-point precision (e.g., dot product)





## **Anytime Automaton**

#### More details in paper:

- Asynchronous/synchronous pipelining
- Data locality with sampling
- Approximate storage techniques
- Thread scheduling

# Evaluation – Methodology

#### **Experiments:**

- IBM Power 780 system
  - 4 POWER7+ cores
  - 32 total hardware threads

#### **Applications:**

- PERFECT and AxBench suites
  - 2D convolution (output sampling, reduced precision<sup>+</sup>, SRAM bit upsets<sup>+</sup>)
  - debayer (output sampling)
  - discrete wavelet transform (loop perforation)
  - histogram equalization (input and output sampling)
  - k-means clustering (output sampling)























The Edward S. Rogers Sr. Department of Electrical & Computer Engineering UNIVERSITY OF TORONTO

141



The Edward S. Rogers Sr. Department of Electrical & Computer Engineering UNIVERSITY OF TORONTO

142





### Evaluation – Discrete Wavelet Transform










## **Evaluation – Summary**



### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.



### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.



#### application execution

### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.



### We propose the **Anytime Automaton**:

> A new computation model for approximate computing.



#### application execution

# Thank you

### The Anytime Automaton

Joshua San Miguel Natalie Enright Jerger



Special thanks to IBM collaborators: Viji Srinivasan, Ravi Nair, Dan Prener