

# Parallel HMMs

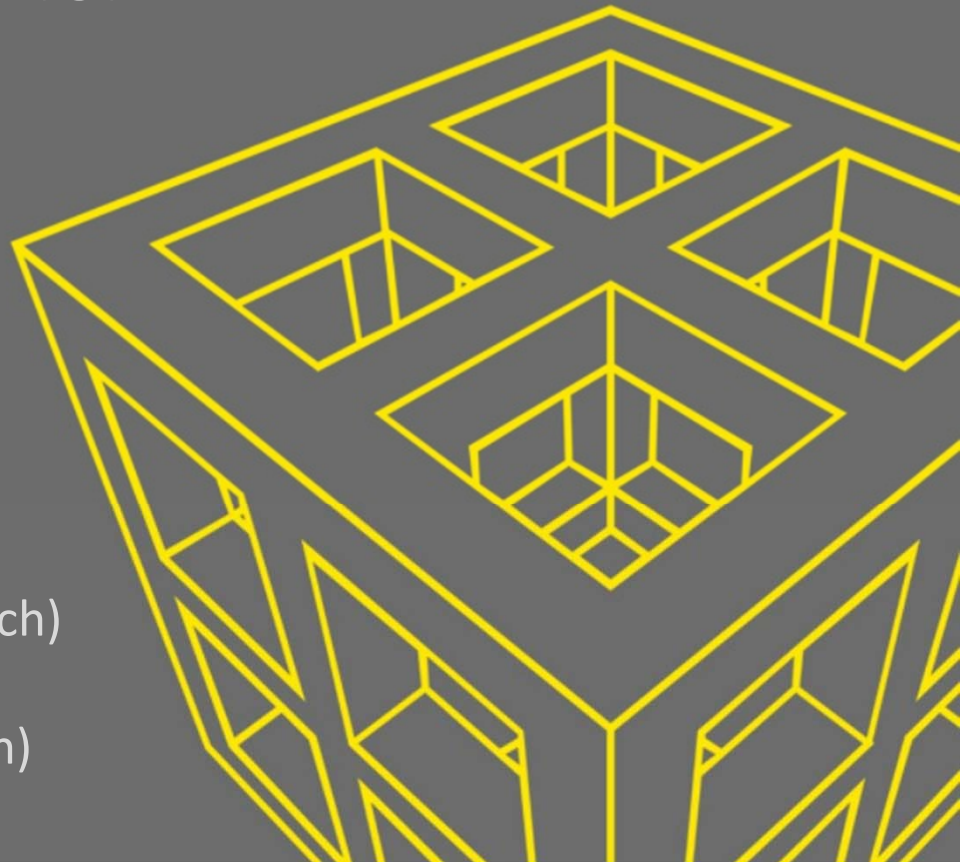
Parallel Implementation of  
Hidden Markov Models for  
Wireless Applications

## Authors

Shawn Hymel (Wireless@VT, Virginia Tech)

Ihsan Akbar (Harris Corporation)

Jeffrey Reed (Wireless@VT, Virginia Tech)

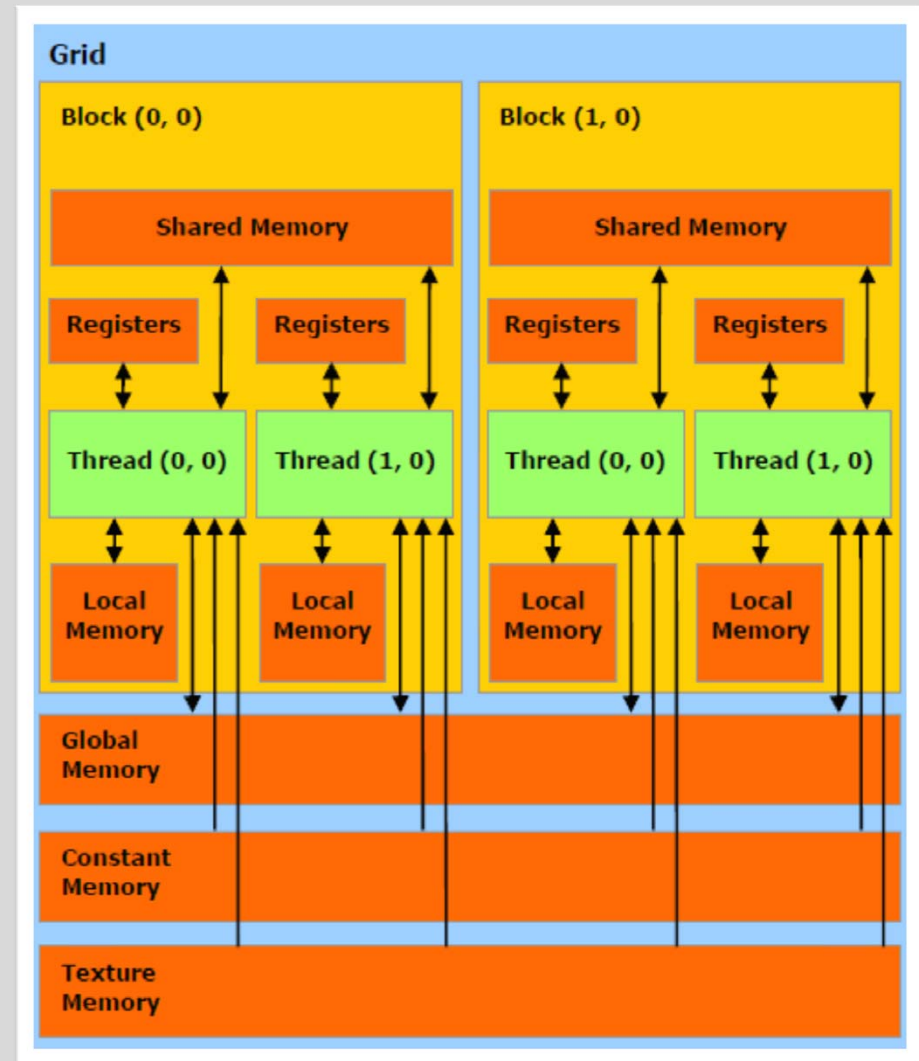


# Agenda

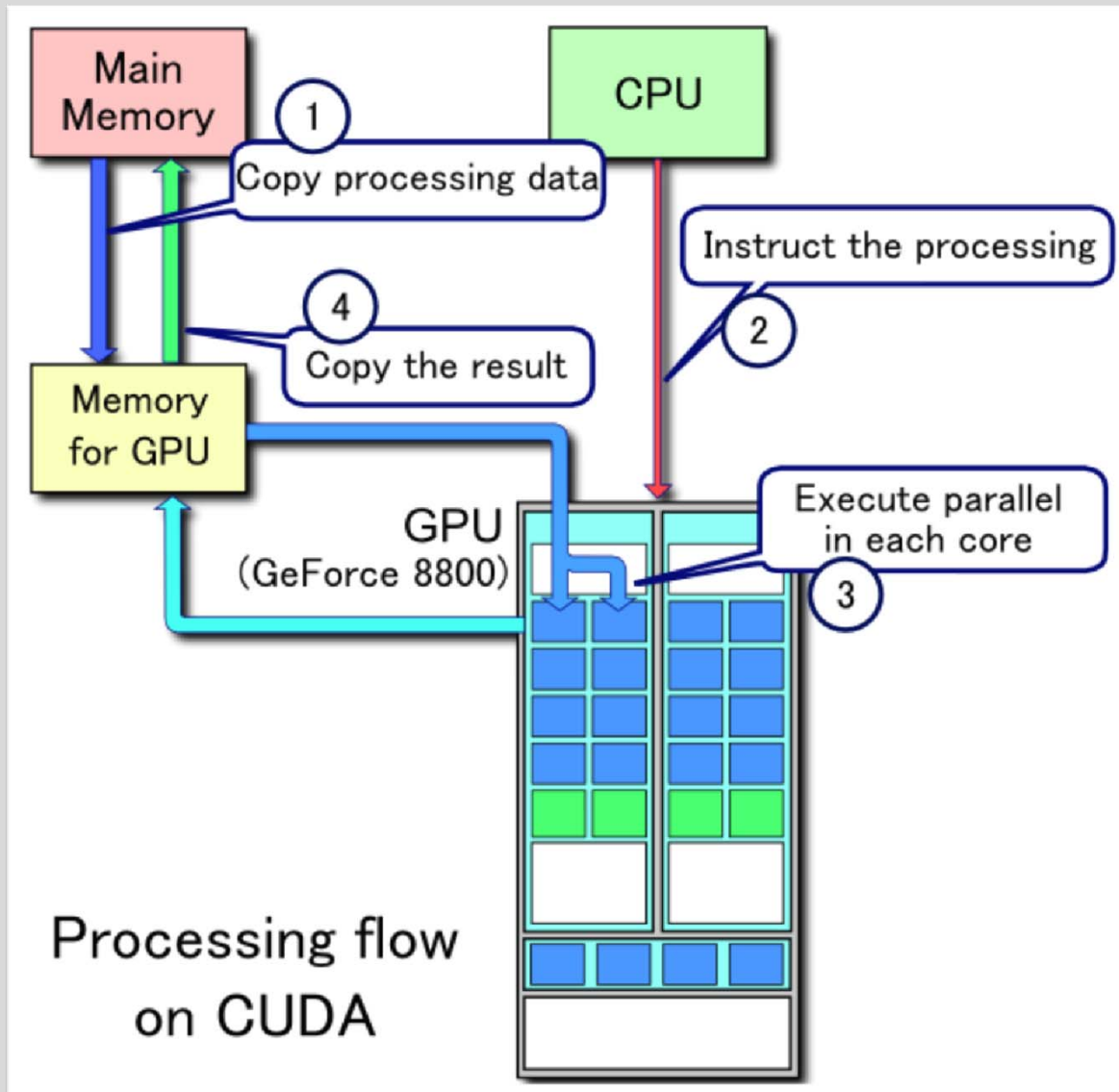
- Overview of GPGPU
- Overview of HMMs
- Parallelization
- Results
- Applications
- Why Is This Useful?

# General-Purpose Processing on GPUs

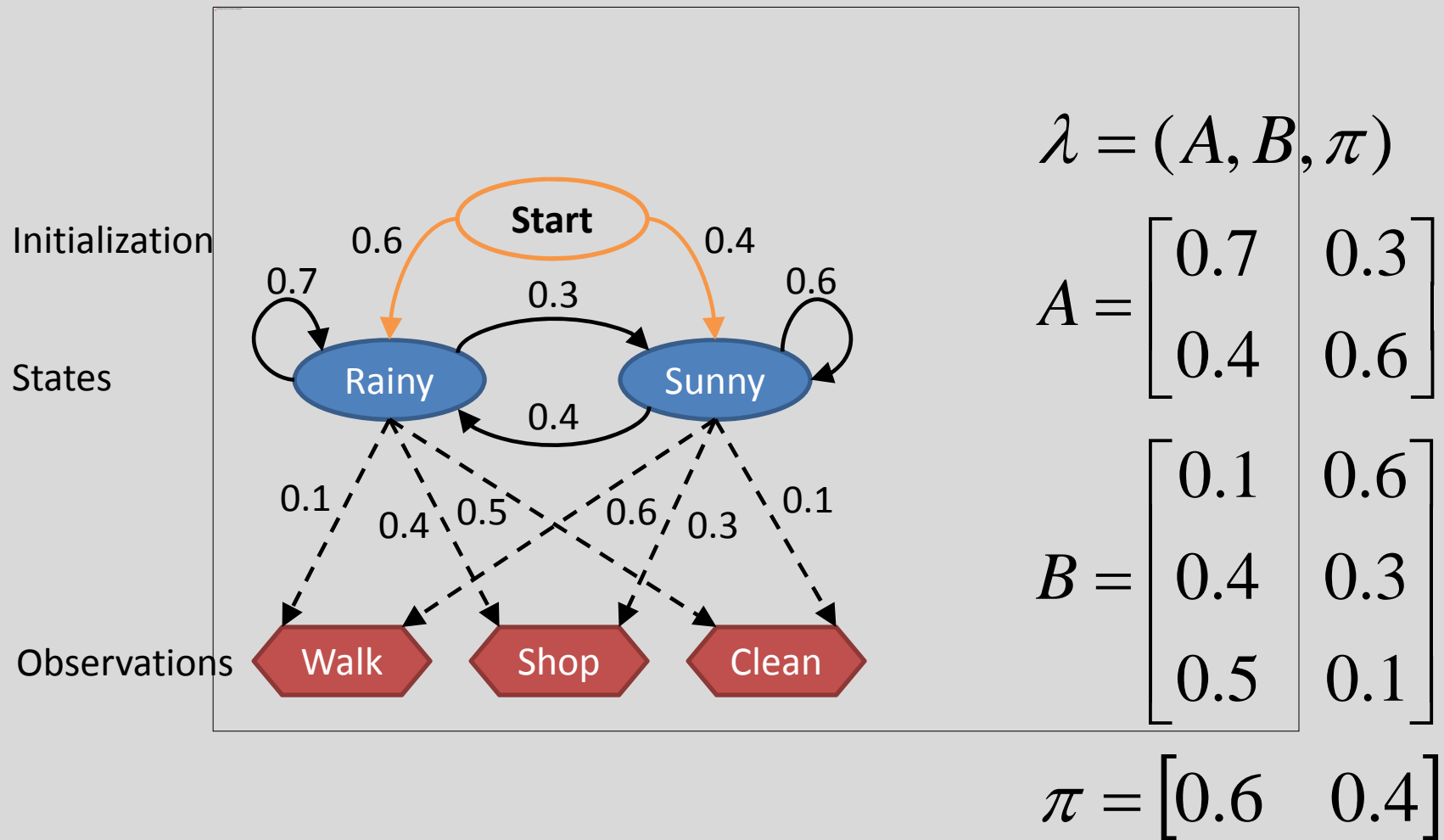
- CUDA-specific
- Important Terms:
  - Threads
  - Blocks
  - Grid



# CUDA Code Flow



# Hidden Markov Model



# HMM Canonical Problems

- Evaluation:  $P(O|\lambda)$ 
  - Forward Algorithm
  - Backward Algorithm
- Find the most likely state sequence
  - Viterbi Algorithm
- Training ( maximize  $P(O|\lambda)$  )
  - Baum-Welch Algorithm

# Forward Algorithm

Given a model and an observation sequence, calculate  $P(O|\lambda)$

- $T$  = number of observations
- $N$  = number of states
- $M$  = number of possible symbols

Initiation:

$$\alpha_1(i) = \pi_i b_i(O_1), i = 1, 2, \dots, N$$

Induction:

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1})$$

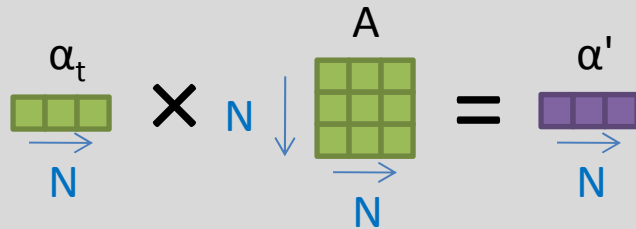
Termination

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$

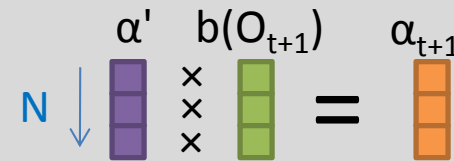
# Example of Parallelization

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1})$$

For all j, matrix multiplication



For all j, element-by-element multiplication



We can perform this step in parallel!

$$O(TN^2) \rightarrow O(T \log N)$$



# Computational Complexity

	Serial	Parallel
<b>Forward Algorithm</b>	$O(TN^2)$	$O(T \log N)$
<b>Viterbi Algorithm</b>	$O(TN^2)$	$O(T \log N)$
<b>Baum-Welch Algorithm</b>	$O(TN^2)$ or $O(TMN)$	$O(T \log N)$

# Test Procedures

- Time execution of each algorithm (C vs. CUDA)
  - Vary states
  - Vary symbols
  - Vary sequence length
- Calculate total energy consumption (C vs. CUDA)
  - PowerTOP software

## Test Hardware

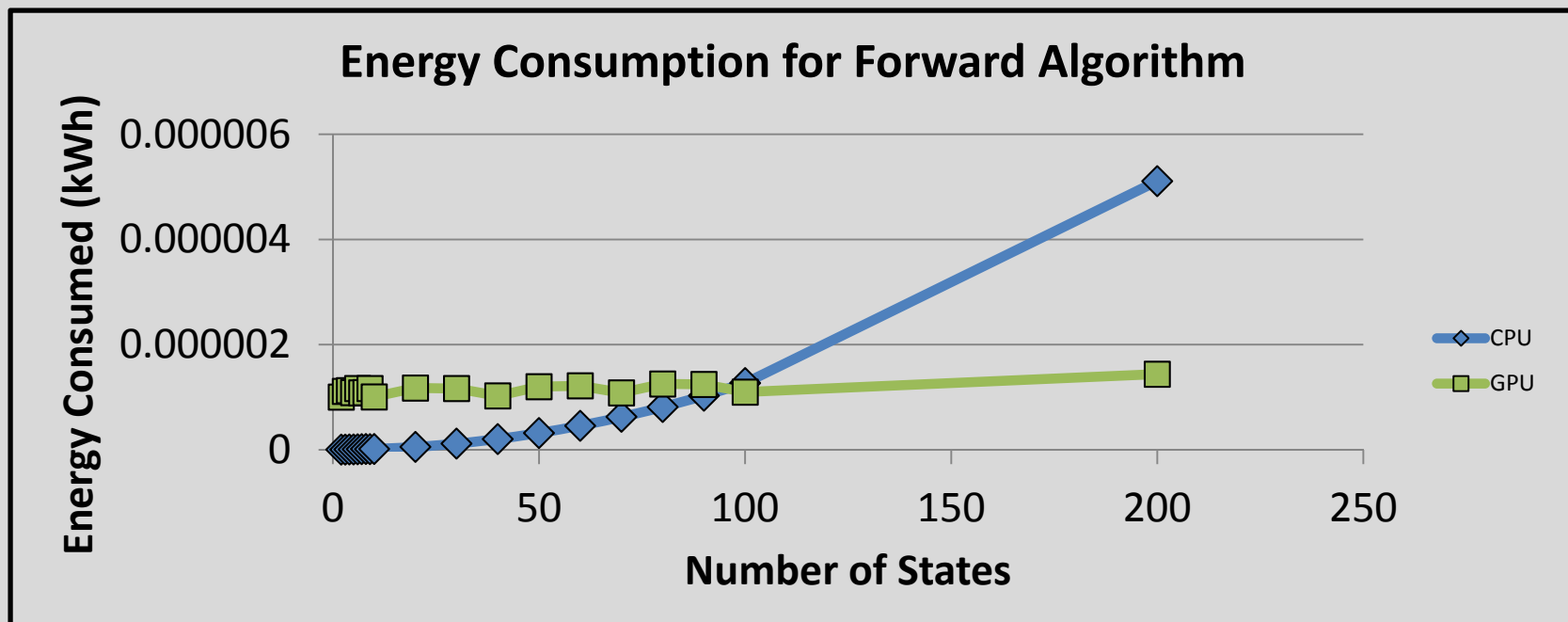
Component	Specification
CPU	Intel Core 2 Duo U7300 @ 1.30GHz
GPU	NVIDIA GeForce GT 335M
GPU Core Speed	450 MHz
GPU Shader Speed	1080 MHz
GPU Memory Speed	1066 MHz
CUDA Cores	72

# Speed Results

Number of States	CPU Runtime (s)	GPU Runtime (s)	Speed Increase
<i>Forward Algorithm</i>			
4	0.001	0.1531	<b>0.007x</b>
40	0.04	0.1393	<b>0.287x</b>
400	4.2816	0.2379	<b>17.99x</b>
4000	534.2028	2.9495	<b>181.12 x</b>
<i>Viterbi Algorithm</i>			
4	0.0033	0.1605	<b>0.021x</b>
40	0.0436	0.1801	<b>0.242x</b>
400	4.2684	1.6595	<b>2.57x</b>
4000	534.5543	116.2531	<b>4.60 x</b>
<i>Baum-Welch Algorithm</i>			
4	0.0021	0.4142	<b>0.005x</b>
40	0.1946	0.4299	<b>0.453x</b>
400	17.6719	0.7502	<b>23.56x</b>
4000	1834.672	28.1271	<b>65.23 x</b>

# Energy Consumption

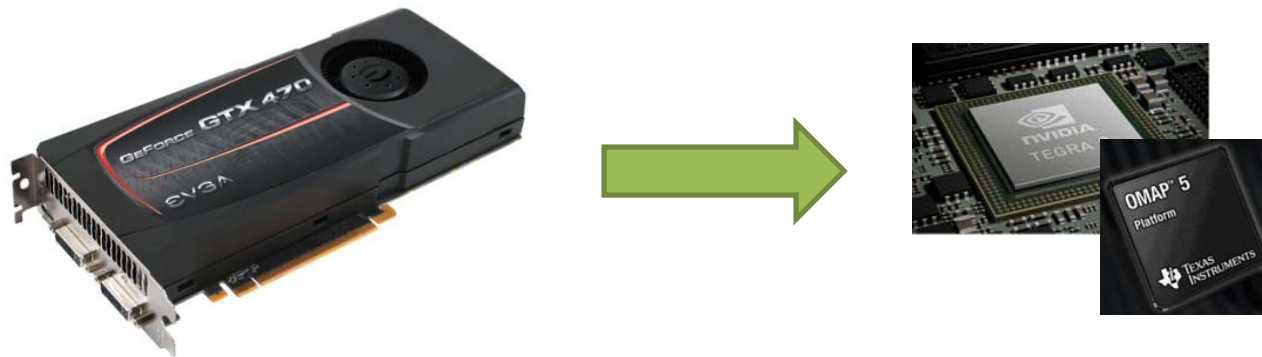
Algorithm	Power (W)		States to Break Even
	C	CUDA	
Forward	18.5	26.5	~100
Viterbi	18.5	29.1	~120
BWA	18.3	26.1	~70



# Applications

- Pattern Recognition
  - Spectrum Sensing
  - Signal Classification
  - Specific Emitter Identification
  - Geolocation
- Modeling
  - Channel Fading
  - Call Drop Prediction

# Why Is This Useful?



- Evolution of GPUs and multi-core processors
  - Smart phones, tablets, SDR
  - Co-processor
- Utilize existing hardware for HMM applications
  - Large number of states
  - 2D/3D HMMs
- Uses in other fields (speech recognition, computer vision)
- Extrapolation to other algorithms (pattern recognition)

# Questions?

## Contact Information

Email: [hymelsr@vt.edu](mailto:hymelsr@vt.edu)

Blog: <http://sgmustadio.wordpress.com/>

Code: <http://code.google.com/p/hmm-cuda/>

## Other Good Resources

cuHMM: <http://code.google.com/p/chmm/>

MATLAB: <http://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>

HTK: <http://htk.eng.cam.ac.uk/>

# Supporting Slide: Reductions

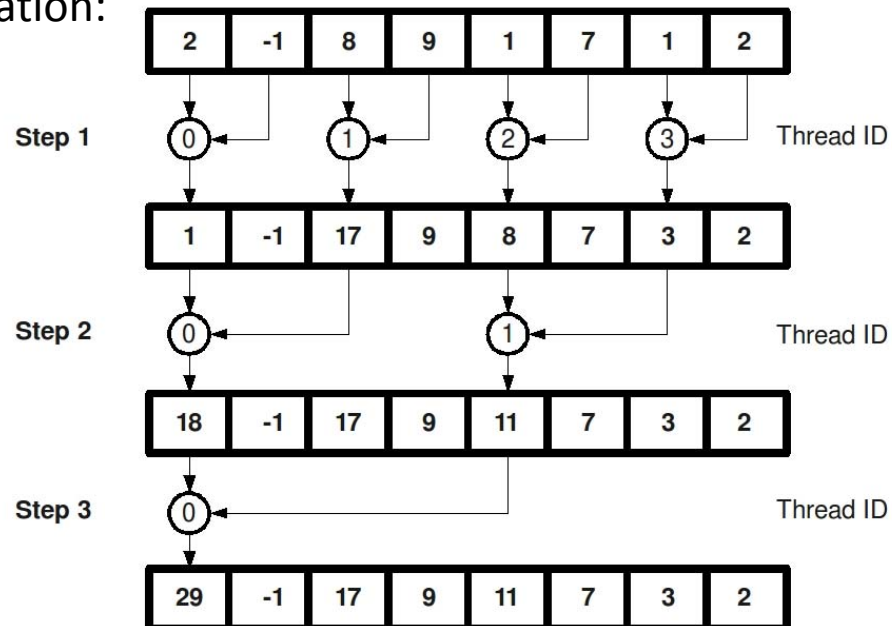
MATLAB example:

```
>> sum(A)
```

C Implementation:

```
sum = 0;  
for (i = 0; i < length; i++) {  
    sum = sum + A[i];  
}
```

Parallelization:

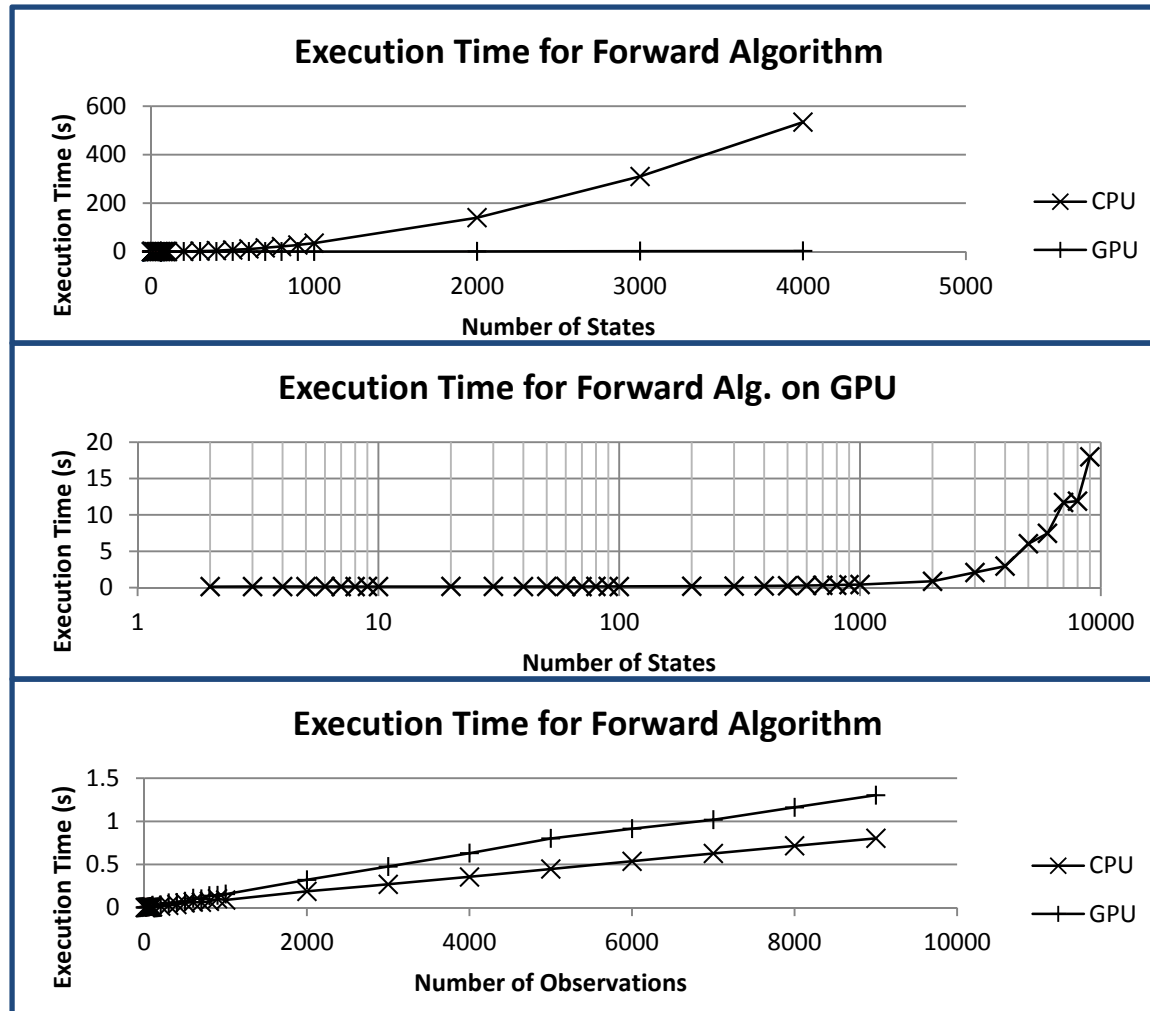


Reducing arrays to a single value (e.g. sum) go from  $O(N)$  to  $O(\log N)$



# Supporting Slide: Timing Results (Forward)

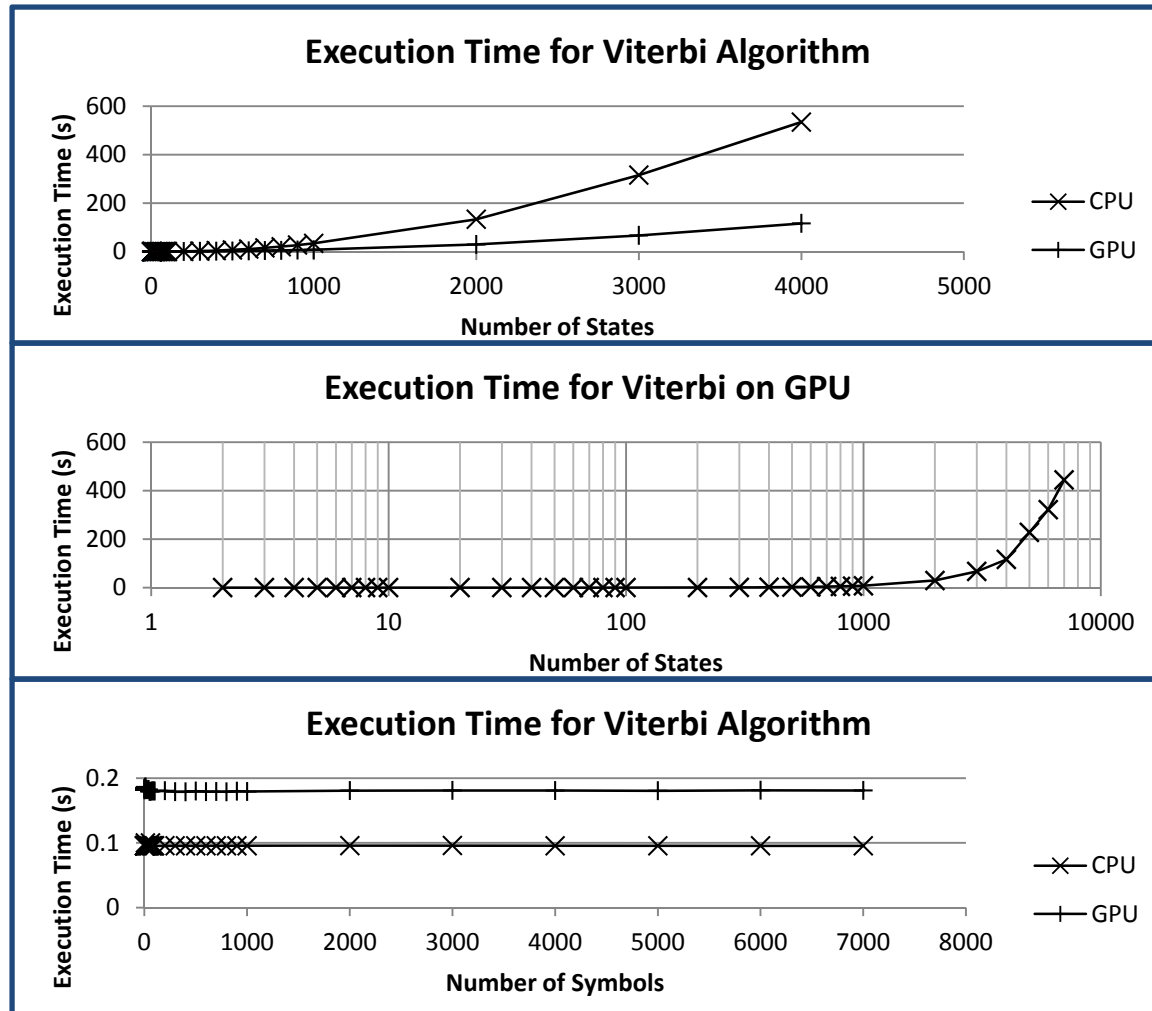
Vary States



Vary Symbols

# Supporting Slide: Timing Results (Viterbi)

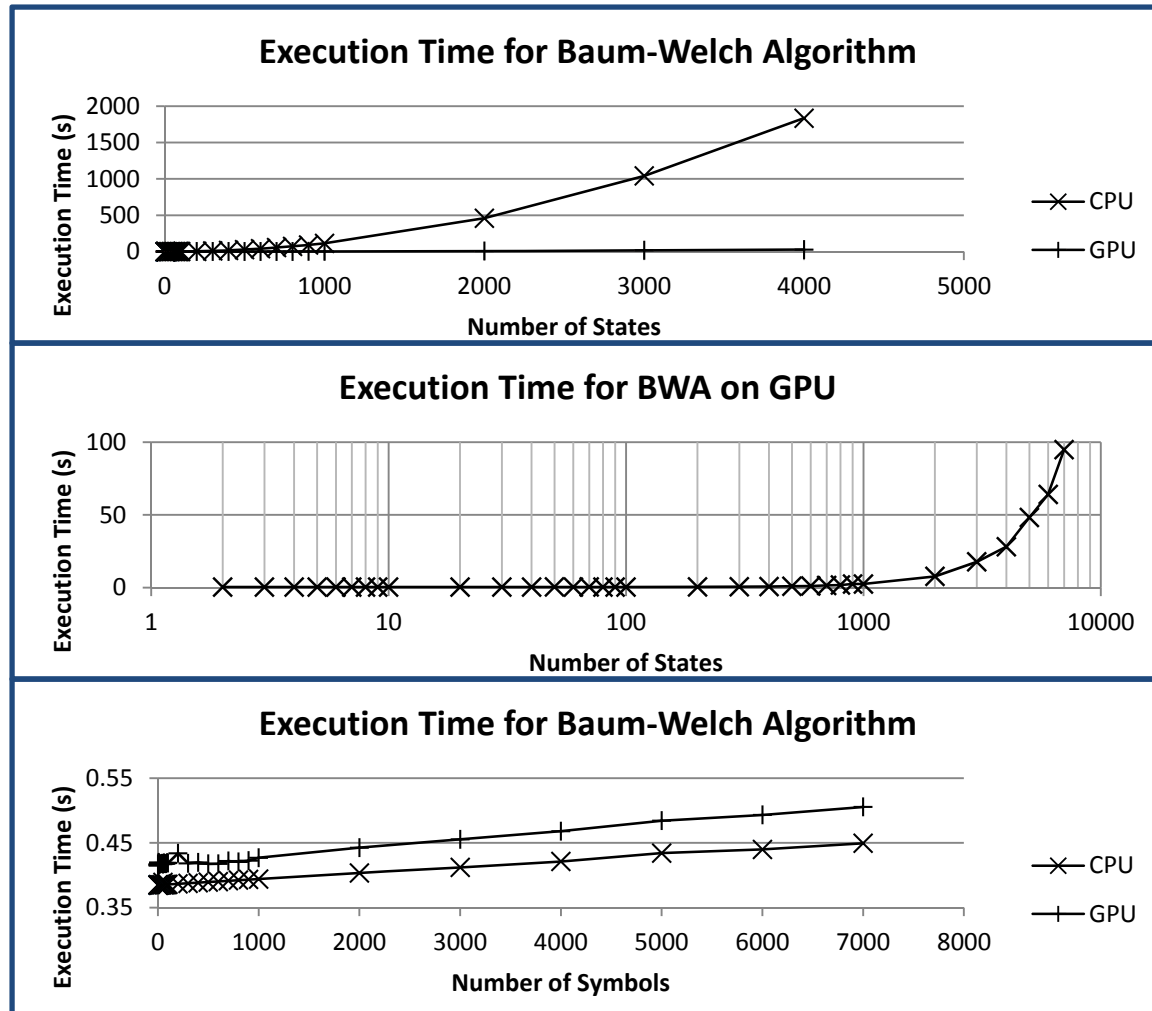
Vary States



Vary Symbols

# Supporting Slide: Timing Results (BWA)

Vary States



Vary Symbols