

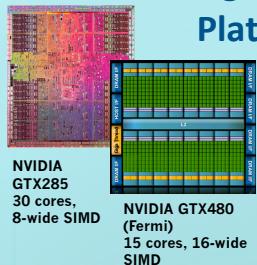


# Exploring Recognition Network Representations for Efficient Speech Inference on Highly Parallel Platforms



Jike Chong, Ekaterina Gonina, Kisun You, Kurt Keutzer, Department of Electrical Engineering and Computer Science, University of California, Berkeley

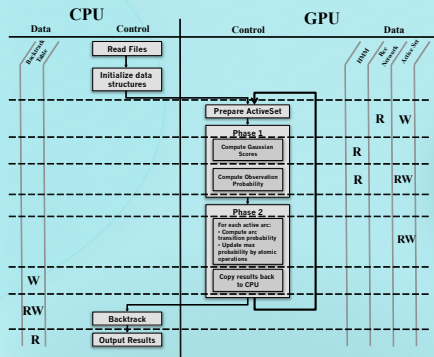
## Maturing Highly Parallel Platforms



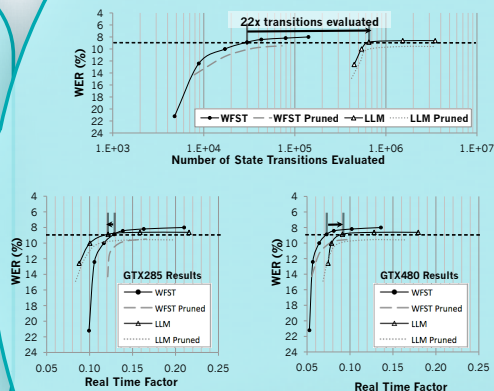
- Architecture trend:
  - Increasing vector unit width
  - Increasing numbers of cores per die
- Maturing HW architecture:
  - Including caches as well as local stores that benefit irregular accesses

Ongoing work investigates performance of alternative approaches to speech recognition on these highly parallel platforms

## Implementation Architecture



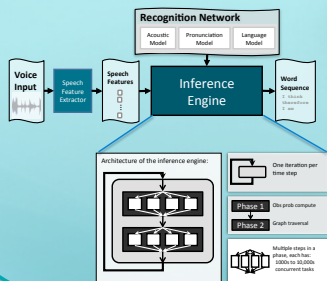
## Evaluation of the Recognition Network Representations



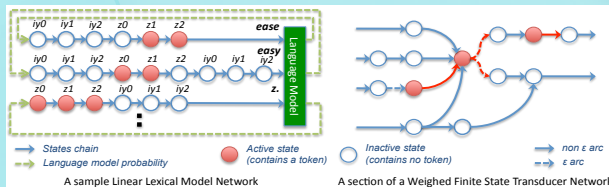
- To achieve the same accuracy:
  - LLM traverses **22x** more state transitions than WFSST
  - On GTX285, LLM is faster
  - On GTX480, WFSST is faster
- Looking at detailed timing:
  - LLM takes 3-5x more time in Graph Traversal, but evaluates 22x more transitions
  - 51% of the execution time in WFSST is spent in gathering data from its irregular data structure

## Speech Recognition Inference Engine Characteristics

- Parallel graph traversal through Recognition network
  - Guided by a sequence of input audio vectors
  - Computing on continuously changing data working set
- Implementation challenges
  - Define a scalable software architecture to expose fine-grained application concurrency
  - Efficiently synchronize between an increasing number of concurrent tasks
  - Effectively utilize the SIMD-level parallelism



## Two Recognition Network Representations



- LLM Network**
  - Chain of triphone states for each pronunciation
  - Each chain constructed using a separate copy of triphone states – many duplications
  - Evaluate possibility of transition from one word to all other words at the end of each triphone chain
- WFSST Network**
  - FSM of composed pronunciation and language models
  - Across-word transitions explicitly represented
  - Encapsulates large amount of information with little redundancy
  - Fewer tokens required to be maintained for target accuracy

	LLM Pruned	LLM	WFSST Pruned	WFSST
# States	123,246	123,246	1,091,295	3,925,931
# Arcs	537,608	1,596,884	2,955,145	11,394,956

- Wall Street Journal 1 Corpus**
  - Based on a 5,000 word vocabulary, 1,350,392 bigrams (291,116 pruned)
  - 3000 16-mixture acoustic models, 39 dim features based on 13 dim MFCC
  - WFSST network is an HCLG model compiled and optimized offline

Real Time Factor	LLM - GTX285	WFSST - GTX285	LLM - GTX480	WFSST - GTX480
8.22x faster than real time	Data Gathering: 18%	Observation Prob: 14%	Graph Traversal: 45%	Sequential Overhead: 23%
7.17x faster than real time	Data Gathering: 51%	Observation Prob: 10%	Graph Traversal: 12%	Sequential Overhead: 26%
10.95x faster than real time	Data Gathering: 17%	Observation Prob: 20%	Graph Traversal: 47%	Sequential Overhead: 16%
11.75x faster than real time	Data Gathering: 51%	Observation Prob: 15%	Graph Traversal: 13%	Sequential Overhead: 19%
Execution speed measured at 8.90% WER				

## Conclusions

- Simpler LLM network representation performs competitively with highly optimized WFSST representation
- WFSST representation is a more concise representation requiring traversal of 1/22 number of state transitions to achieve the same accuracy
- Per state transition LLM gathers data 53-65x faster and evaluates transitions 4.7-6.4x faster than WFSST
- Uncoalesced memory accesses are still a major bottleneck in implementations using the WFSST representation

Emergence of highly parallel platforms brings forth an opportunity to reevaluate computational efficiency of speech recognition approaches.

Thanks to Nelson Morgan, Andreas Stolcke, and Adam Janin at ICSI for insightful discussions and continued support in the infrastructure used in this research.

This research is supported in part by an Intel Ph.D. Fellowship.

This research also supported in part by Microsoft (Award #024263) and Intel (Award #024894) funding and by matching funding by U.C. Discovery (Award #DIG07-10227).