Electrical & Computer
**ENGINEERING**

2007 STARC Forum

# Towards Speech Recognition in Silicon:
## The Carnegie Mellon *In Silico Vox* Project

Rob A. Rutenbar

Professor, Electrical & Computer Engineering

rutenbar@ece.cmu.edu

© R.A. Rutenbar 2007

**Carnegie Mellon**

---

**CarnegieMellon**

# Speech Recognition Today

- **Quality = *OK*     Vocab = *large***

- **Quality = *poor*    Vocab = *small***

- **Commonality: all software apps**

© Rob A. Rutenbar 2007                                           Slide 2

---

**CarnegieMellon**

# Today's Best *Software* Speech Recognizers

- **Best-quality recognition is computationally *hard***
  - ❯ For speaker-independent, large-vocabulary, continuous speech

- **1-10-100-1000 rule**
  - ❯ For **~1X** real-time recognition rate
  - ❯ For **~10%** word error rate (90% accuracy)
  - ❯ Need **~100 MB** memory footprint
  - ❯ Need **~100 W** power
  - ❯ Need **~1000 MHz** CPU

- **This proves to be very *limiting* …**
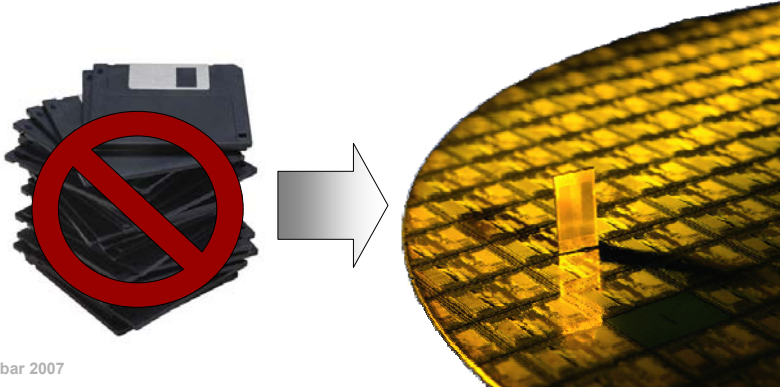
© Rob A. Rutenbar 2007                                                      Slide 3

**CarnegieMellon**

# The Carnegie Mellon *In Silico Vox* Project

- **The thesis: It's time to liberate speech recognition from the unreasonable limitations of software**

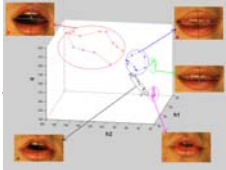- **The solution: *Speech recognition in silicon***



© Rob A. Rutenbar 2007

**CarnegieMellon**

## Aside: About the Name "*In Silico Vox*"

- **■** *In Vivo*
  - ◤ Latin: an experiment done in a living organism

- **■** *In Vitro*
  - ◤ Latin: an experiment done in an artificial lab environment

- **■** *In Silico*
  - ◤ (*Not* real Latin): an experiment done via computation only

- **■** *Vox*
  - ◤ Latin: voice, or word

© Rob A. Rutenbar 2007    Slide 5



**CarnegieMellon**

## About This Talk

- **■ Some philosophy**
  - ◤ Why silicon? Why now? Why us (CMU)?

- **■ A quick tour: How speech recognition works**
  - ◤ What happens in a recognizer

- **■ An SoC architecture**
  - ◤ Stripping away all CPU stuff we don't need, focus on essentials

- **■ Results**
  - ◤ ASIC version: Simulation results
  - ◤ FPGA version: Live, running hardware-based recognizer

© Rob A. Rutenbar 2007    Slide 6

**CarnegieMellon**

## About This Talk

- **Some philosophy**
  - Why silicon?   Why now?   Why us (CMU)?

- **A quick tour:  How speech recognition works**
  - What happens in a recognizer

- **An SoC architecture**
  - Stripping away all CPU stuff we don't need, focus on essentials

- **Results**
  - ASIC version:   Simulation results
  - FPGA version:   Live, running hardware-based recognizer

© Rob A. Rutenbar 2007                                                                 Slide 7

---

**CarnegieMellon**

## Why Silicon?   Why Now?

**Why?   Two reasons:**

- **History**
  - We have some successful **historical** examples of this migration

- **Performance**
  - Tomorrow's compelling apps need **100X – 1000X** more performance
  - (Not going to happen in software)

© Rob A. Rutenbar 2007                                                                 Slide 8

CarnegieMellon

## History:    Graphics Engines

- **Nobody paints pixels in software anymore!**
  - Too limiting in max performance.   Too inefficient in power.

True on the desktop (& laptop)

NVIDIA® GeForce® 7950 GX2

http://www.nvidia.com

…and on your cellphone too

MV8602

http://www.mtekvision.com

© Rob A. Rutenbar 2007                    Slide 9



CarnegieMellon

## Performance: Next-Gen Compelling Applications

**Audio-mining**
- **Very fast recognizers – much faster than realtime**
- **App:  search large media streams (DVD) quickly**

FIND: **"Hasta la vista, baby!"**

**Hands-free appliances**
- **Very portable recognizers – high quality result on << 1 watt**
- **App:  interfaces to small devices, cellphone dictation**

"send email to arnold – let's do lunch…"

© Rob A. Rutenbar 2007                    Slide 10

---

**CarnegieMellon**

## Silicon Solution:  Speed *and* Power Wins

- **A famous graph from Prof. Bob Brodersen of Berkeley**
  - Study looked at 20 designs published at ISSCC, from 1997-2002
  - In slightly older technologies, relative to today:  180nm – 250nm
  - Dedicated designs up to **10,000X better** energy efficiency (MOPS/mW)



© Rob A. Rutenbar 2007                                                      Slide 11

---

**CarnegieMellon**

## Recent Example: Parallel Radio Baseband DSP

- **90nm CMOS: adaptive DSP for multipath MIMO channel**
  - Power efficiency = **2.1GOPS/mW**
  - Area efficiency = **20GOPS/mm$^2$**

| Technology | 90nm CMOS |
|---|---|
| Core area | $1.9 \times 1.9$ mm |
| Die area | $2.3 \times 2.3$ mm |
| Pad count | 120 |
| IO/core $V_{DD}$ | 1V / 0.4V |
| Cell count | 420,304 |
| Frequency | 100 MHz |
| P (act/leak) | 30mW / 4mW |
| Efficiency | 2.1GOPS/mW |



1st path, $\alpha_1 = 1$

Tx array

Rx array

2nd path, $\alpha_2 = 0.6$

x                                   y

Data rate up to 250Mbps over 16 sub-carriers
Measured 34mW @ VDD=385mV

(Source:  Prof. Dejan Markovitz, UCLA)

© Rob A. Rutenbar 2007                                                      Slide 12

---

**Rob A. Rutenbar**
**Carnegie Mellon University**

CarnegieMellon

## About This Talk

- **Some philosophy**
  - Why silicon?  Why now?  Why us (CMU)?

- **A quick tour:  How speech recognition works**
  - What happens in a recognizer

- **An SoC architecture**
  - Stripping away all CPU stuff we don't need, focus on essentials

- **Results**
  - ASIC version:  Simulation results
  - FPGA version:  Live, running hardware-based recognizer

© Rob A. Rutenbar 2007

Slide 15



CarnegieMellon

## How Speech Recognition Works

Adaptation to environment/speaker

Adaptation

ADC

Sampling

Filter1
Filter2
Filter3
FilterN
DSP

$x_1$
$x_2$
$x_3$
$x_n$

"ao"

Acoustic

Word
...
Rob  R AO B
Bob  B AO B
...

Language

Rob → says

HMM / Viterbi Search
Language Model Search

"Rob"

Feature extraction

Feature vector

Feature Scoring

Acoustic units → Words → Language

**Acoustic Frontend**   **Scoring**   **Backend Search**

© Rob A. Rutenbar 2007

Slide 16

**(3) Search: Speech Models are *Layered* Models**

Language X Words X Acoustic → Layered Search

YES

NO

words

/Y/ → /EH/ → /S/

/N/ → /OW/

"acoustic units"

"sub-acoustic units"

Power
Spectrum
Waveform
Pitch

1 frame of sampled sound

Classical methods (HMMs, Viterbi) and idiosyncracies

© Rob A. Rutenbar 2007    Slide 19



**Context Matters:    At Bottom -- *Triphones***

- **English has ~50 atomic sounds (phones) but we recognize ~50x50x50 context-dependent triphones**
  - Because "I" sound in "five" is different than the "I" in "nine"

Five    $F(-, I)_{cross-word}$    $I(F, V)_{word-internal}$    $V(I, -)_{cross-word}$

Nine    $N(-, I)_{cross-word}$    $I(N, N)_{word-internal}$    $N(I, -)_{cross-word}$

**"I" in "five"  ≠  "I" in  "nine"**

© Rob A. Rutenbar 2007    Slide 20

---

**CarnegieMellon**

## Similar for Other Languages, like Japanese

■ ...but different basic building blocks (different **phones**)

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ア | a Modifier | カ ガ | ka ga | サ ザ | sa za | タ ダ | ta da |
| ア | a | キ ギ | ki gi | シ ジ | shi ji | チ ヂ | chi q |
| イ | i Modifier | ク | ku | ス ズ | su zu | ツ ヅ | tsu |
| イ | i | グ | gu | セ ゼ | se ze | テ デ | te de |
| ウ | u Modifier | ケ ゲ | ke ge | ソ ゾ | so zo | ト ド | to do |
| ウ | u | コ ゴ | ko go | | | | |
| エ | e Modifier | | go | | | | |
| エ | e | | | | | | |
| オ | o Modifier | | | | | | |
| オ | o | | | | | | |

マ ミ ム メ モ — ma mi mu me mo

| ハ バ パ | ha ba pa | ラ リ ル レ ロ | ra ri ru re ro | ワ ヲ ン ヴ | wa wo nn vu |
|---|---|---|---|---|---|
| ヒ ビ ピ | hi bi pi | ャ ヤ | ya Modifier / ya | カ ケ ヴ ギ ゞ ヂ | ka ke va vi ve vo |
| フ ブ プ | hu bu pu | ュ ユ | yu Modifier / yu | | |
| ヘ ベ ペ | he be pe | ョ ヨ | yo Modifier / yo | | |
| ホ ボ ポ | ho bo po | | | | |

. Middle dot
— Prolonged sound
ヽ Iteration
ヾ Voiced iteration

Source: Microsoft Speech API 5.3
Japanese Phonemes
http://msdn2.microsoft.com/en-us/library/ms720568.aspx

© Rob A. Rutenbar 2007

Slide 21

---

**CarnegieMellon**

## Example of Different Phone Building Blocks

■ **Let's use my name as an example**
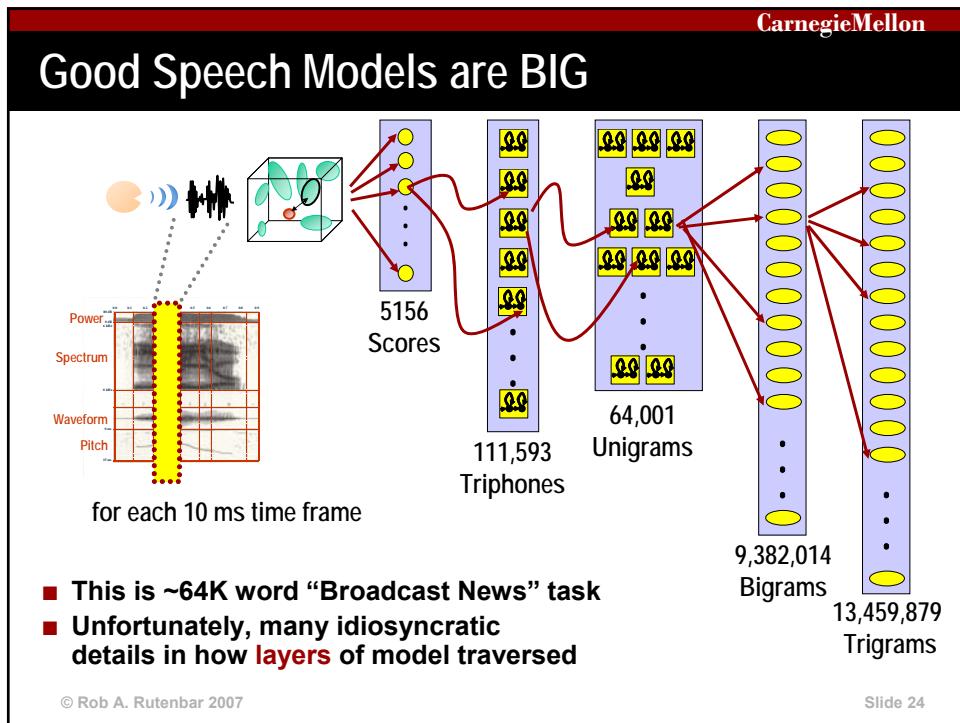
■ **English:**    /r/    /OO/    /t/    /n/    /b/    /ar/

■ **Japanese:**    /ル/    /ー/    /テ/    /ン/    /バ/    /ー/

■ **Aside:**
  ◥ Japanese has some reputation as being an "easier" language for automatic recognition
  ◥ Mapping from basic sounds (mora) to words is simpler than English

© Rob A. Rutenbar 2007

Slide 22

---

Also Context at Top: *N-gram* Language Model

Suppose we have vocabulary
{ W1, W2, W3, W4, …}

Lets us calculate likelihood of word **W3** *after* **W2** *after* **W1**

© Rob A. Rutenbar 2007 — Slide 23



Good Speech Models are BIG

5156 Scores

111,593 Triphones

64,001 Unigrams

9,382,014 Bigrams

13,459,879 Trigrams

for each 10 ms time frame

- **This is ~64K word "Broadcast News" task**
- **Unfortunately, many idiosyncratic details in how layers of model traversed**

© Rob A. Rutenbar 2007 — Slide 24

---

**CarnegieMellon**

## Where Does *Software* Spend its Time?

- **CPU time for CMU Sphinx 3.0**
  - Prior studies targeted less capable versions (v1, v2)
  - Tools: SimpleScalar & Intel Vtune
  - 64K-word "Broadcast News" benchmark

- **So: It's all *backend***

6% 15% 19% 27% 33%

~0% of time!

Adaptation to environment/speaker

Acoustic | Word | Language
Rob R AO B
Bob B AO B
Rob → says

"ao"

ADC — DSP — Feature extraction — Feature vector — Feature Scoring — HMM / Viterbi Search Language Model Search — "Rob"

Sampling

Acoustic units → Words → Language

© Rob A. Rutenbar 2007   Slide 25

---

**CarnegieMellon**

## Memory Usage? SPHINX 3.0 vs Spec CPU2000

- **Cache sizes**
  - L1: 64 KB, direct mapped
  - DL1: 64 KB, direct mapped
  - UL2: 512 KB, 4-way set assoc

- **So…**
  - Terrible locality (no surprise, graph search + huge datasets)
  - Load dominated (no surprise, reads a lot, computes a little)
  - Not an insignificant footprint

|  | SPHINX 3.0 | Gcc | Gzip | Equake |
|---|---|---|---|---|
| Cycles | 53 B | 55B | 15 B | 23 B |
| IPC | 0.69 | 0.29 | 1.05 | 0.7 |
| **Instruction Mixes** | | | | |
| Loads | 0.27 | 0.25 | 0.2 | 0.27 |
| Stores | 0.05 | 0.15 | 0.09 | 0.08 |
| Branch's | 0.14 | 0.2 | 0.17 | 0.12 |
| **Branch Misprediction Rates** | | | | |
|  | 0.025 | 0.07 | 0.08 | 0.02 |
| **Cache Miss Rates** | | | | |
| DL1 | 0.04 | 0.02 | 0.02 | 0.03 |
| L2 | 0.48 | 0.06 | 0.03 | 0.30 |
| **Memory Footprint** | | | | |
|  | 64 MB | 24 MB | 186 MB | 42 MB |

© Rob A. Rutenbar 2007   Slide 26

---

## About This Talk

- **Some philosophy**
  - Why silicon?   Why now?   Why us (CMU)?

- **A quick tour:  How speech recognition works**
  - What happens in a recognizer

- **An SoC architecture**
  - Stripping away all CPU stuff we don't need, focus on essentials

- **Results**
  - ASIC version:   Simulation results
  - FPGA version:   Live, running hardware-based recognizer

© Rob A. Rutenbar 2007      Slide 27

---

## This Talk:   How to Get to *Fast*...

**Audio-mining**
- **Very fast recognizers – much faster than realtime**
- **App:  search large media streams (DVD) quickly**

FIND: **"Hasta la vista, baby!"**
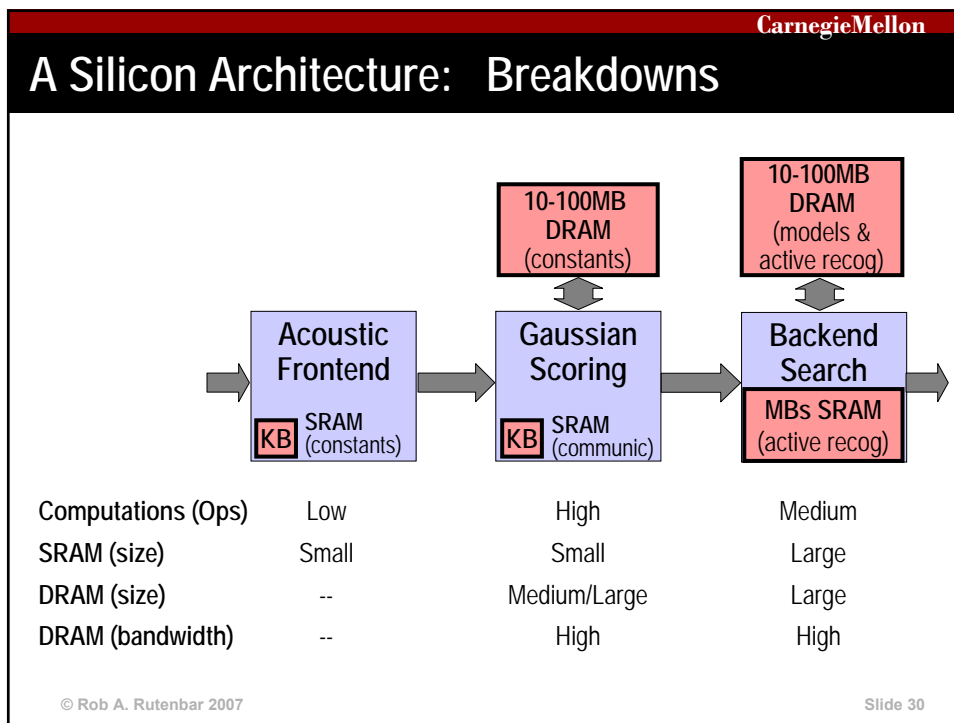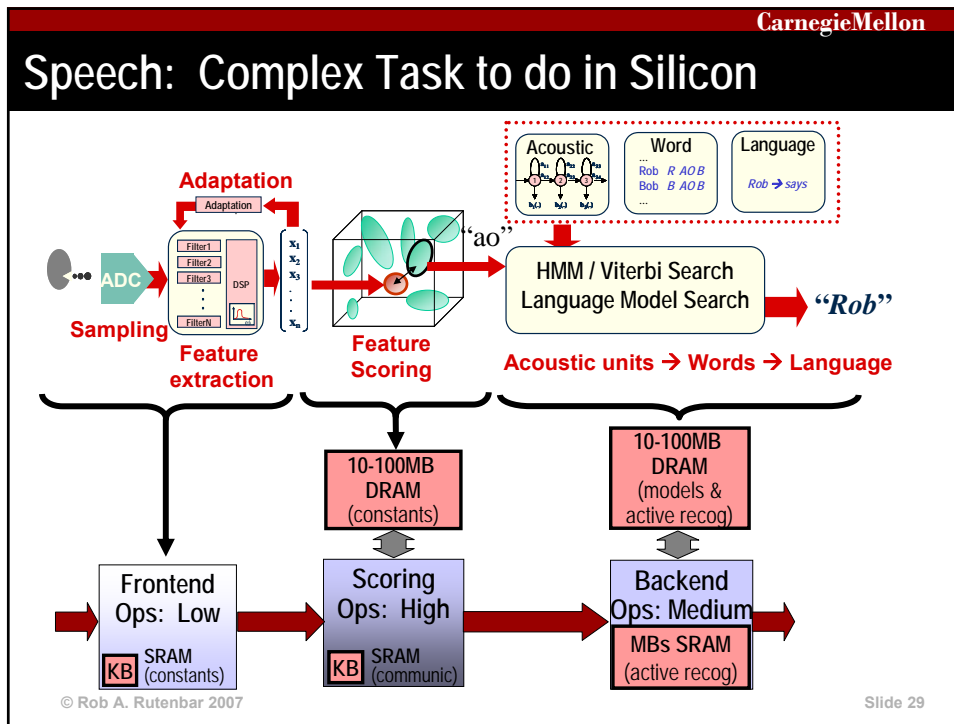
**Hands-free appliances**
- **Very portable recognizers – high quality result on << 1 watt**
- **App:  interfaces to small devices, cellphone dictation**

"send email to arnold – let's do lunch..."

© Rob A. Rutenbar 2007      Slide 28

---

Speech: Complex Task to do in Silicon



A Silicon Architecture: Breakdowns

|  | Acoustic Frontend | Gaussian Scoring | Backend Search |
|---|---|---|---|
| Computations (Ops) | Low | High | Medium |
| SRAM (size) | Small | Small | Large |
| DRAM (size) | -- | Medium/Large | Large |
| DRAM (bandwidth) | -- | High | High |

**CarnegieMellon**

# Essential Implementation Ideas

- **Custom precision, everywhere**
  - Every bit counts, no extras, no floating point – all fixed point

- **(Almost) no caching**
  - Like graphics chips: fetch from SDRAM, do careful data placement
  - (Little bit of caching for bandwidth filtering on big language models)

- **Aggressive pipelining**
  - If we can possibly overlap computations – we try to do so

- **Algorithm transformation**
  - Some software computations are just bad news for hardware
  - Substitute some "deep computation" with hardware-friendly versions

Slide 31

**CarnegieMellon**

# Example: Aggressive Pipelining

**Pipelined *Get-HMM/Viterbi* and *Transition* stages**

| Fetch Word | Fetch HMM/Viterbi | Transition/Prune/Writeback | | time |
| | | Fetch HMM/Viterbi | Transition/Prune/Writeback | Language Model |

- Word #1
- Word #2
- HMM #1
- HMM #2
- HMM #3

**Pipelined *Get-Word* and *Get-HMM* stages**

| Fetch Word | Fetch HMM/Viterbi | Transition/Prune/Writeback | |
| | Fetch Word | Fetch HMM/Viterbi | Transition/Prune/Writeback | Language Model |

**Pipelined *non-LanguageModel* and *LanguageModel* stages**

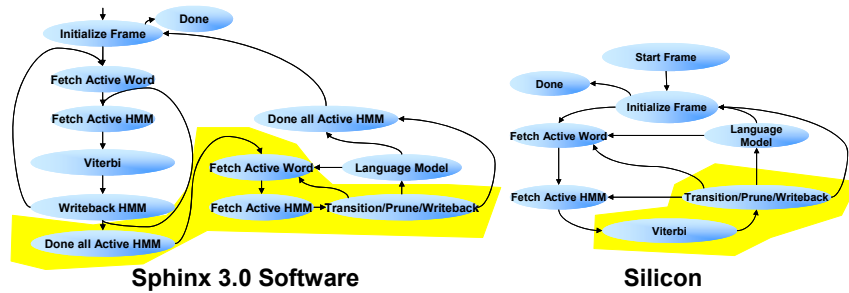| Fetch Word | Fetch HMM/Viterbi | Transition/Prune/Writeback | Language Model |
| | Fetch Word | Fetch HMM/Viterbi | Transition/Prune/Writeback |
| | | Fetch HMM/Viterbi | Transition/Prune/Writeback |

Slide 32

**Example: Algorithmic Changes**

- **Acoustic-level pruning threshold**
  - **Software**: Use best score of *current* frame (after Viterbi on Active HMMs)
  - **Silicon**: Use best score of *previous* frame (nixes big temporal bottleneck)
- **Tradeoffs**
  - Less memory bandwidth, can pipeline, little pessimistic on scores

**Sphinx 3.0 Software**          **Silicon**

© Rob A. Rutenbar 2007          Slide 33

---

**About This Talk**

- **Some philosophy**
  - Why silicon? Why now? Why us (CMU)?

- **A quick tour: How speech recognition works**
  - What happens in a recognizer

- **An SoC architecture**
  - Stripping away all CPU stuff we don't need, focus on essentials

- **Results**
  - ASIC version:   Simulation results
  - FPGA version:   Live, running hardware-based recognizer

© Rob A. Rutenbar 2007          Slide 34

---

## Design Flow: C++ Cycle Simulator → Verilog

- **2006 benchmark: 5K-word "Wall Street Journal" task**
- **Cycle sim results:**
  - No accuracy loss; not quite 2X @ 125MHz ASIC clock
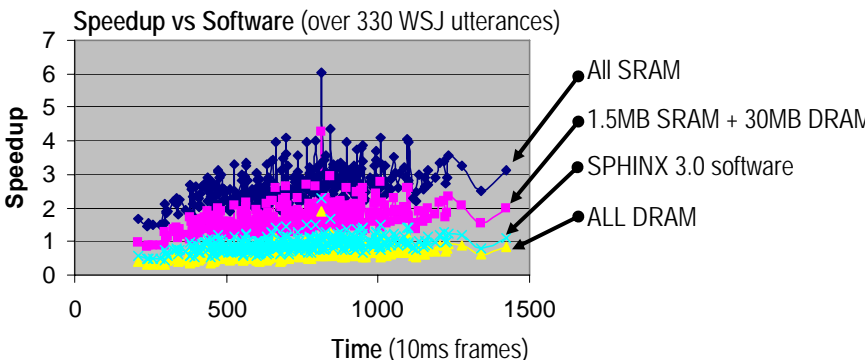  - Backend search needs: ~1.5MB SRAM, ~30MB DRAM

| Recognizer Engine | Word Error Rate (%) | Clock (GHz) | Speedup Over Real Time (bigger is better) |
|---|---|---|---|
| Software: Sphinx 3.3 (fast decoder) | 7.32% | 1 GHz | 0.74X |
| Software: Sphinx 4 (single CPU) | 6.97% | 1 GHz | 0.82X |
| Software: Sphinx 4 (dual CPU) | 6.97% | 1 GHz | 1.05X |
| Software: Sphinx 3.0 (single CPU) | 6.707% | 2.8 GHz | 0.59X |
| Hardware: Our Proposed Recognizer | 6.725% | 0.125 GHz | 1.67X |

© Rob A. Rutenbar 2007    Slide 35

---

CarnegieMellon

## Aside: Bit-Level Verification Hurts (A Lot)

- **We have newfound sympathy for others doing silicon designs that handle large media streams**
  - Generating these sort of tradeoff curves: CPU days → weeks

Speedup vs Software (over 330 WSJ utterances)

All SRAM
1.5MB SRAM + 30MB DRAM
SPHINX 3.0 software
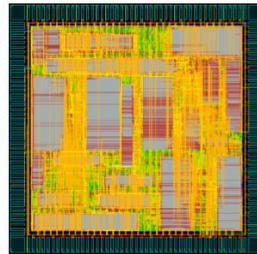ALL DRAM

Speedup / Time (10ms frames)

© Rob A. Rutenbar 2007    Slide 36

Aside: Pieces of Design = Great Class Projects

- CMU student team: Patrick Chiu, David Fu, Mark McCartney, Ajay Panagariya, Chris Thomas
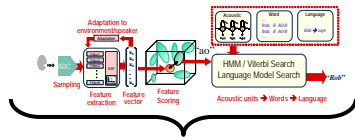
Floorplan   Final Layout   Final Stats

© Rob A. Rutenbar 2007    Slide 37



A Complete Live Recognizer: FPGA Demo

- In any "system design" research, you reach a point where you just want to see it work – *for real*

- Goal: *Full recognizer 1 FPGA + 1 DRAM*

**Xilinx XC2VP30 FPGA**
**[13969 slices / 2448 Kb]**
**Utiliz: 99% of slices**
**45% of Block RAM**
**~3MB DDR DRAM**
**50MHz clk; ~200Mb/s IO**

- A benchmark that fits on chip
  - 1000-word "Resource Mgt" task
  - Slightly simplified: no tri-grams
  - Slower: not real time, ~2.3X slower
  - Resource limited: slices, mem bandwidth

© Rob A. Rutenbar 2007    Slide 38

## FPGA Experimental Results

- **Aside: as far as we know, this is the *most complex* recognizer architecture ever fully mapped into a running, hardware-only form**

    Slide 39

---

CarnegieMellon

## Summary

- **Software is too constraining for speech recognition**
  - Evolution of graphics chips suggests alternative: Do it in silicon
  - Compelling performance and power reasons for silicon speech recog

- **Several "*in silico vox*" architectures in design**
  - SoC and FPGA versions
  - ~10X realtime speedup architecture in progress at CMU

- **Reflections**
  - Some of the most interesting experiences happen when you get people from very different backgrounds – silicon + speech – on same team