# Speech technologies for interactive mobile applications – a primer

Jason D. Williams
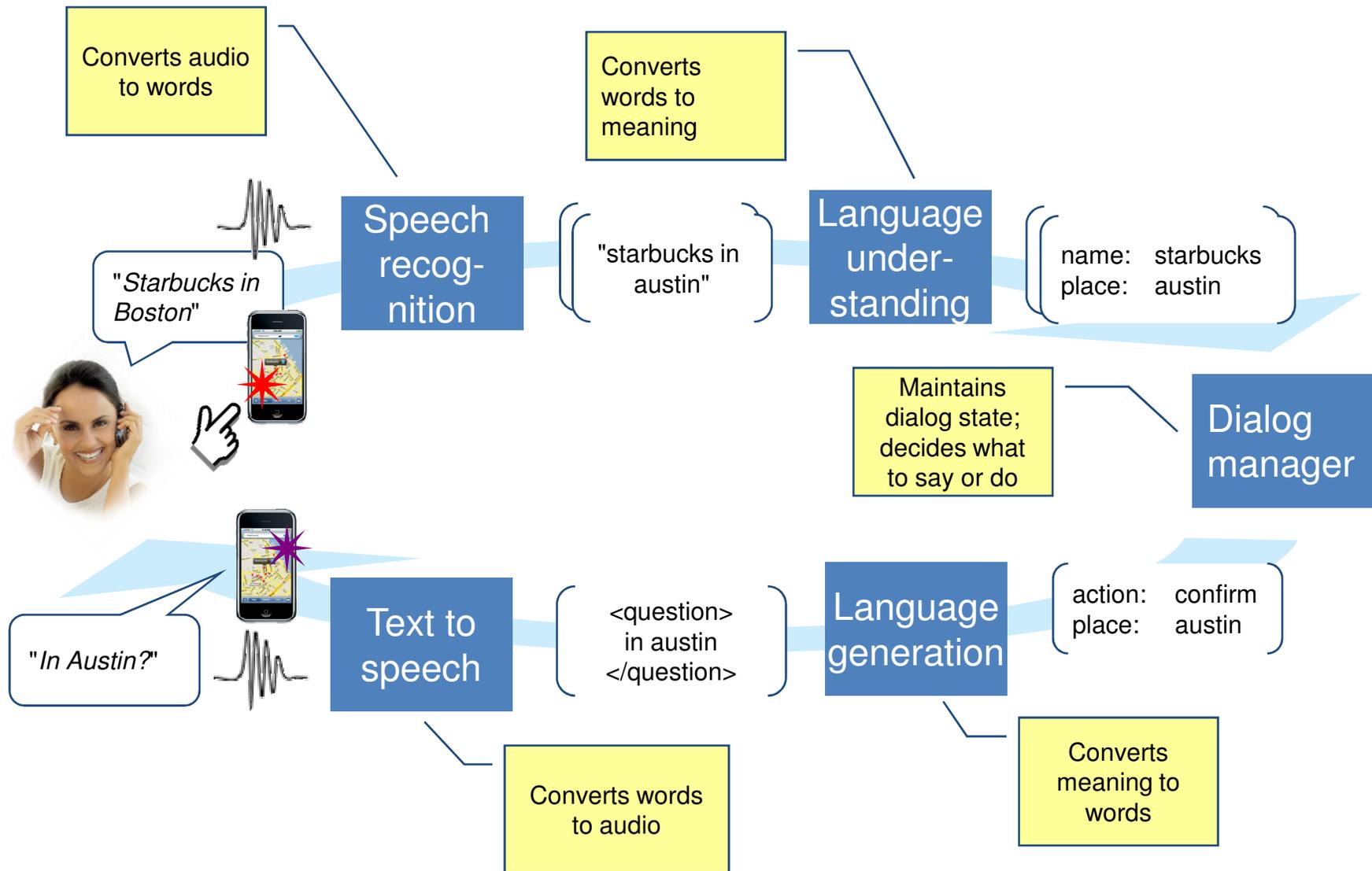

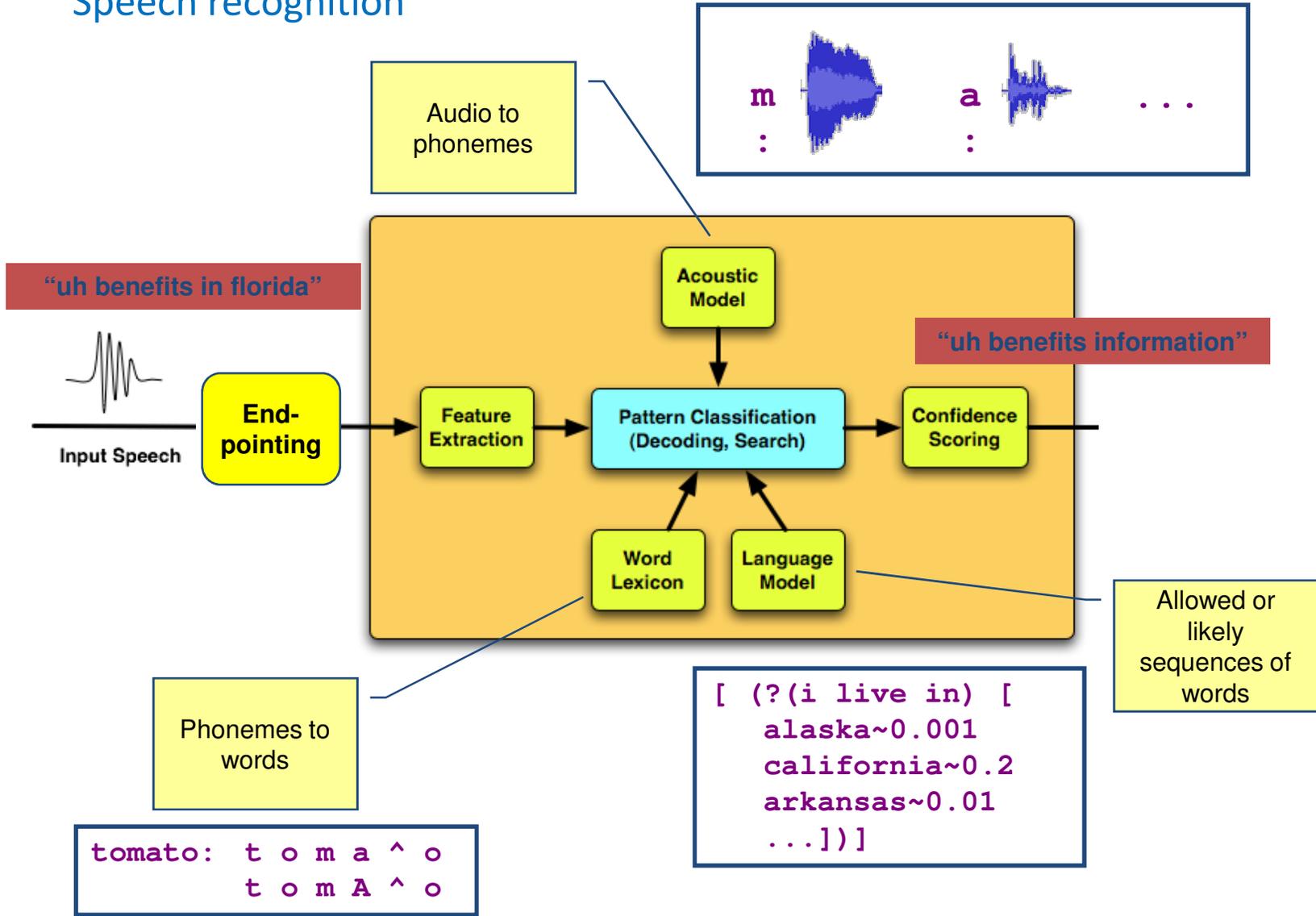at&t

AVIxD / IxD workshop – August 2011

# Roadmap

- A quick tour of the toolbox
- Technology and usage issues (that influence design)
- A few pointers on getting started

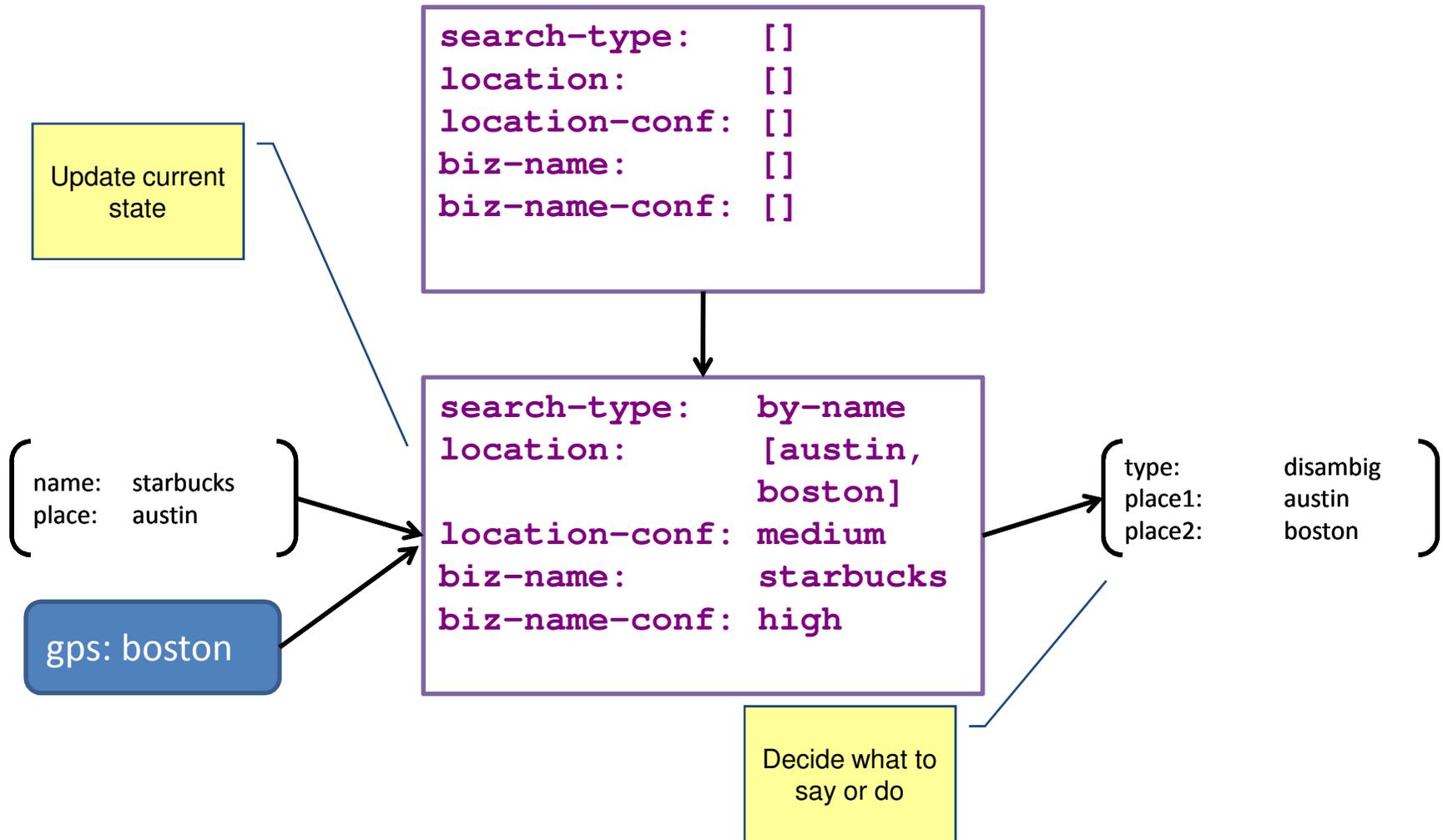# What is a spoken dialog system?

# Speech recognition



Audio to phonemes

m: a: . . .

"uh benefits in florida"

Input Speech

**End-pointing**

Feature Extraction

Acoustic Model

Pattern Classification (Decoding, Search)

Confidence Scoring

"uh benefits information"

Word Lexicon

Language Model

Allowed or likely sequences of words

Phonemes to words

```
tomato:  t o m a ^ o
         t o m A ^ o
```

```
[ (?(i live in) [
    alaska~0.001
    california~0.2
    arkansas~0.01
    ...])]
```

## Language understanding

- Example input/output:

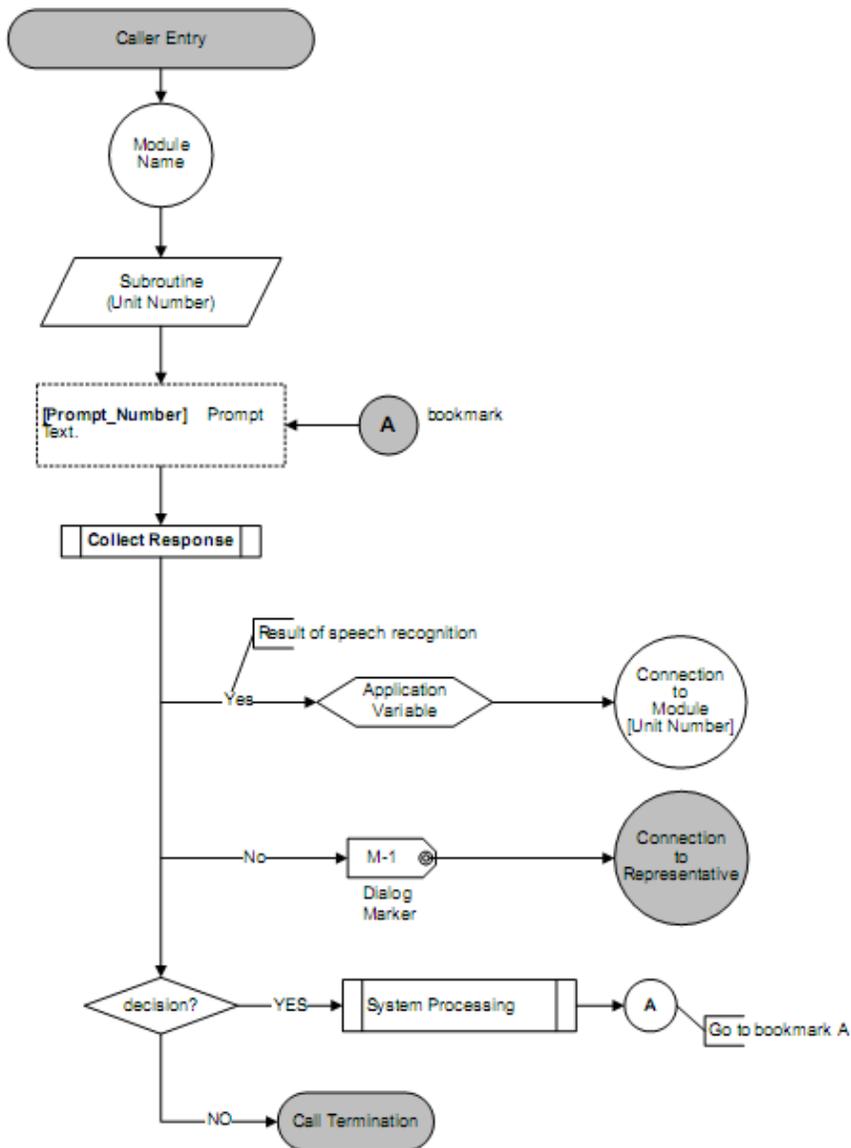| TRANSCRIPTION (input) | INTERPRETATION (output) |
|---|---|
| f_m_l_a claim | <call_request "fmla"> <action "claim"> |
| payroll | <ambig_key "payroll"> |
| customer service | [none] |
| i'd [fragment] ~ | [none] |
| f_m_l_a agent | <call_request "fmla"> <action "operator"> |
| employment | <call_request "employment"> |
| [side_speech] | [none] |
| more examples | <command "moreoptions"> |
| medical benefits | <ambig_key "healthplan"> |

- Can be done with rules, or by a pattern classifier

# Dialog manager

Update current state

```
search-type:    []
location:       []
location-conf:  []
biz-name:       []
biz-name-conf:  []
```

```
name:   starbucks
place:  austin
```

gps: boston

```
search-type:    by-name
location:       [austin,
                boston]
location-conf:  medium
biz-name:       starbucks
biz-name-conf:  high
```

```
type:    disambig
place1:  austin
place2:  boston
```

Decide what to say or do

# How dialog systems are designed today



Typical commercial spoken dialog system contains ~100 pages of flowchart

# Language generation

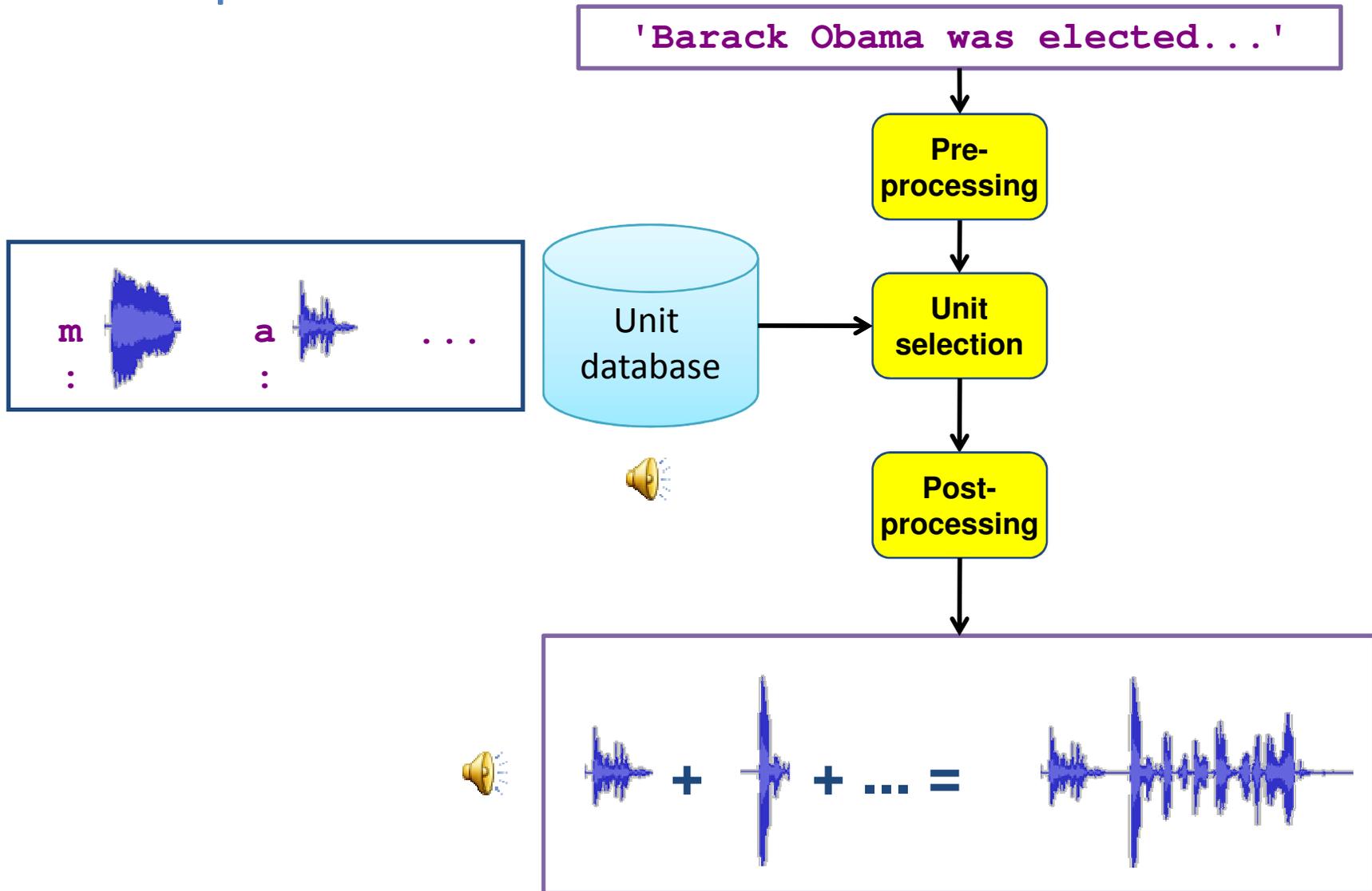$$\left(\begin{array}{ll} \text{type:} & \text{disambig} \\ \text{place1:} & \text{austin} \\ \text{place2:} & \text{boston} \end{array}\right)$$

```python
if (action.type == 'disambig'):
  out.text = 'Sorry, was that %s or %s?' %
(action.place1,action.place2)

elif (action.type == 'confirm'):
  out.text = '%s, is that right?' % (action.item)
  out.wavs = [GetWav(action.item),is_that_right.wav]

elif (action.type == 'display'):
  out.text = '%s in %s' %
(action.business,action.bizname)
  out.display = DisplayMap(action.lat,action.lon)

...
```

**Text-to-speech**

'Barack Obama was elected...'

Pre-processing

m a ...

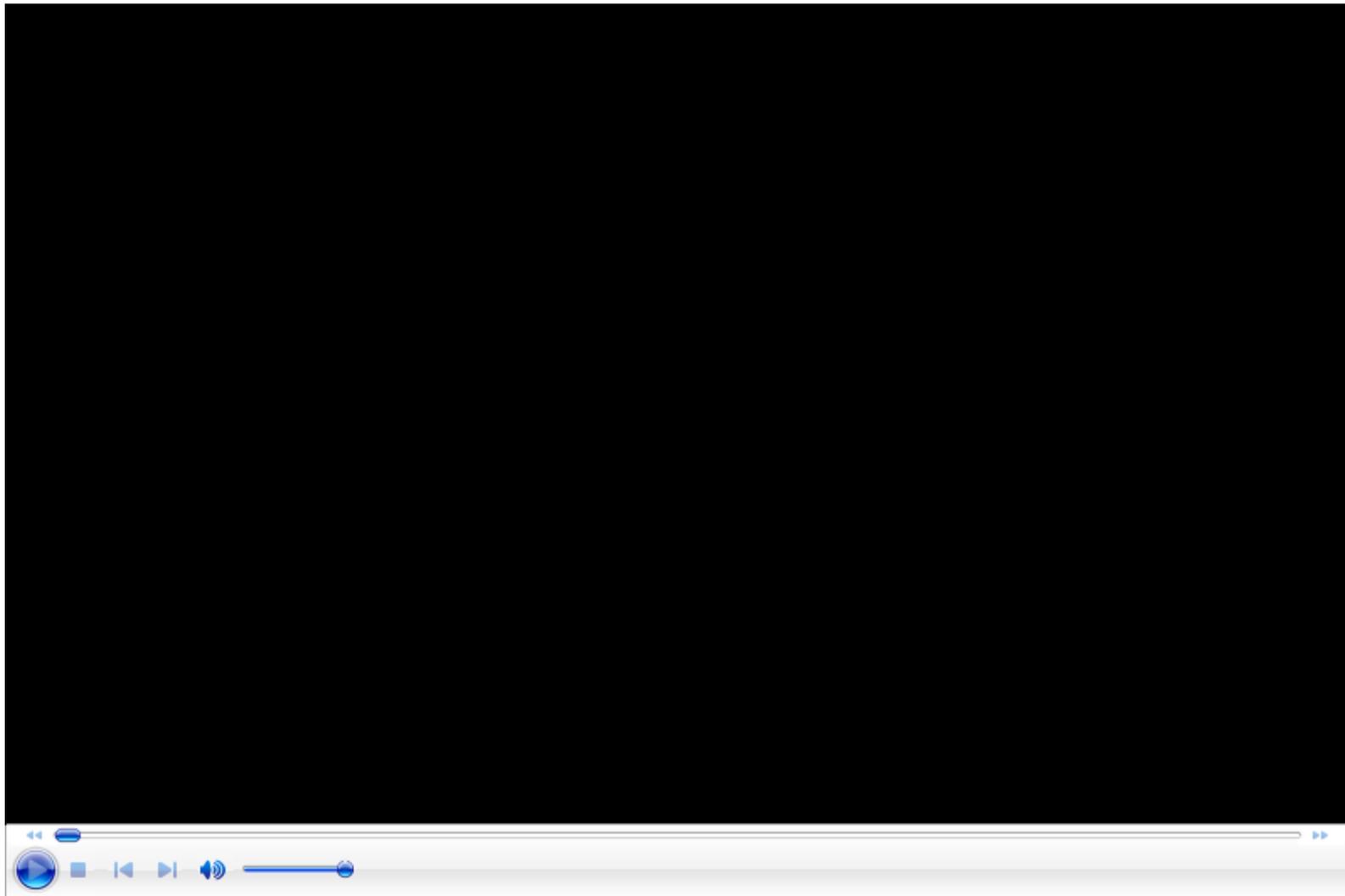Unit database

Unit selection

Post-processing

Putting it all together: example

- Uni-modal telephone-based spoken dialog system

  - SLM speech recognition

  - Language understanding (~250 categories)

  - Speaker verification (~300K users)

# An example multi-modal system

# Two central challenges for deploying speech technology

AVIxD / IxD workshop – August 2011

# ASR/SLU errors are common

| Grammar | Yes/no | City & state | How may I help you? |
|---------|--------|--------------|---------------------|

*Source: Two different deployed commercial applications running two different speech recognizers*

# ASR/SLU errors are common

| Grammar | Yes/no | City & state | How may I help you? |
|---|---|---|---|
| In-grammar/<br>in-domain accuracy | 99.8% | 85.1% | 89.5% |

*Source: Two different deployed commercial applications running two different speech recognizers*

# ASR/SLU errors are common

| Grammar | Yes/no | City & state | How may I help you? |
|---|---|---|---|
| In-grammar/ in-domain accuracy | 99.8% | 85.1% | 89.5% |
| % in-grammar/ in-domain | 92.3% | 91.0% | 86.8% |
| Overall accuracy | 92.1% | 77.6% | 77.7% |

*Source: Two different deployed commercial applications running two different speech recognizers*

# ASR/SLU errors are common

| Grammar | Yes/no | City & state | How may I help you? |
|---|---|---|---|
| In-grammar/ in-domain accuracy | 99.8% | 85.1% | 89.5% |
| % in-grammar/ in-domain | 92.3% | 91.0% | 86.8% |
| Overall accuracy | 92.1% | 77.6% | 77.7% |
| Accepted utts (False accepts) | 89.6% (1.8%) | 60.3% (4.9%) | 73.3% (8.3%) |

*Source: Two different deployed commercial applications running two different speech recognizers*
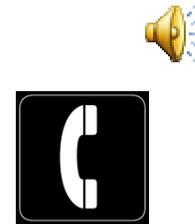
# The "theory of mind" problem

|  | A real human | Graphical user interface | Spoken dialog system |
|---|---|---|---|
| **What can it understand?** | Anything you can explain | Only the buttons you can press | The contents of the grammar |
| **How do I know a task *can't* be performed?** | "I can't do that." | There is no button for it. | It's not in the grammar |

Users must think simultaneously about what language the system can understand, and what the system can do – they must form a "theory of mind" about the dialog system

# Responses to "How may I help you?"

- Silences and hesitations while users think
  - Leads to end-pointing problems
  - Leads to users confusing themselves
- "Robot" language (hence examples, "speak naturally")
  - Example 1
  - Example 2
- Recognition errors confused with competences
  - > "i need to sign up for a get off benefit" *[no parse]*
  - > "i would like to enroll in a get one" *[no parse]*
  - > "i would like to get help with my dental insurance" <HELP>
  - > "dental insurance" <INSURANCE>

# Getting started with speech technology

AVIxD / IxD workshop – August 2011

# How do I add ASR to my mobile application?

- ASR on mobile devices is usually a bad idea
  - ASR will kill your battery
  - Mobile processors are probably too small
  - Can't benefit from cross-user acoustic data
- Probably better to run ASR in the cloud
  - Stream audio to server; server sends back ASR result
  - Example: AT&T Speech Mash-up

# Language models: constrain/weigh word sequences

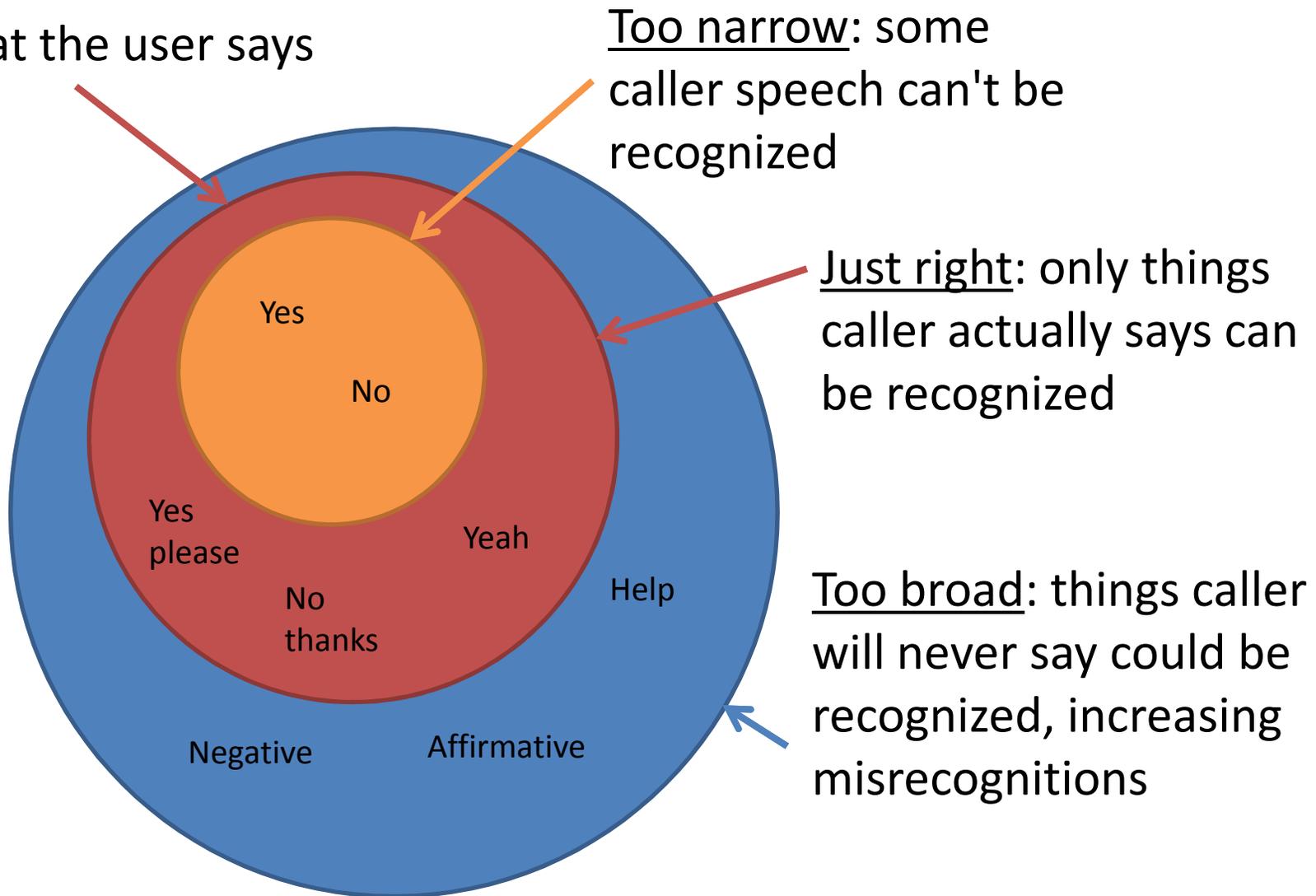Usually, each "dialog state" has its own language model (LM)

Rule-based LMs

```
GET_STATE [
   (?(i live in) [
    alaska~0.001
    california~0.2
    arkansas~0.01
    ...
   ])]
```

Statisical LMs

```
i want 0.0035
want to 0.0023
want benefits 0.0034
want need 0.000002
```

# What to include in the language model?
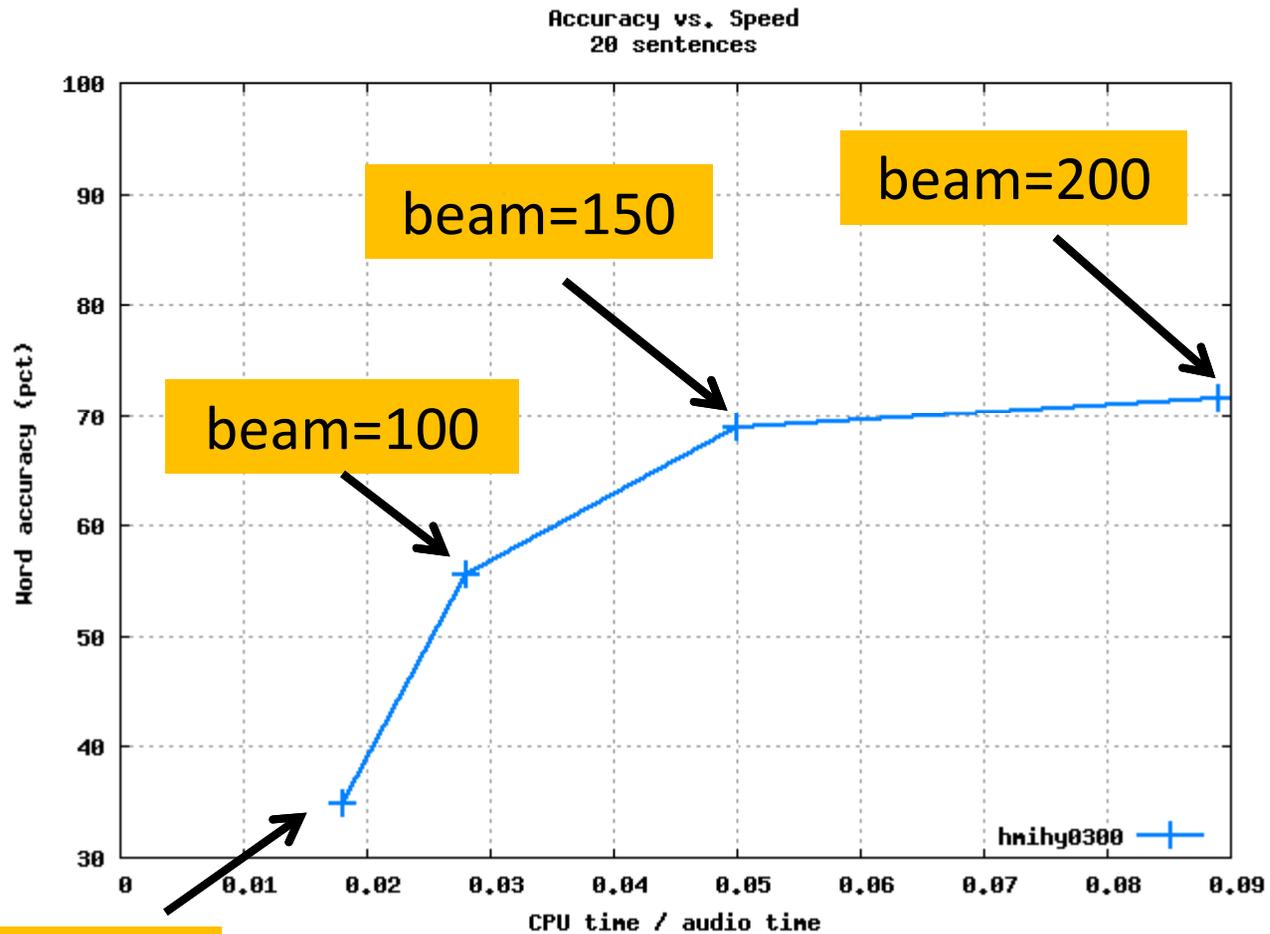
What the user says

Too narrow: some caller speech can't be recognized

Just right: only things caller actually says can be recognized

Too broad: things caller will never say could be recognized, increasing misrecognitions

Yes

No

Yes please

Yeah

No thanks

Help

Negative

Affirmative

# Setting parameters: beam width

# Another parameter: language model scale factor

LOW

HIGH



Too low: good acoustic matches but words make no sense; low accuracy

Just right: a little acoustic information is ignored in favor of words that make sense; best accuracy

Too high: word sequences make sense but too much acoustic information is ignored; low accuracy

Example   "mac store santa monica california"

| LM Scale | Result |
| --- | --- |

Example   "mac store santa monica california"

| LM Scale | Result |
|:---:|:---|
| 0 | mm aq ck storr sade a mon in ck a california |

Example  "mac store santa monica california"

| LM Scale | Result |
|----------|--------|
| 0 | mm aq ck storr sade a mon in ck a california |
| 0.01 | maxtor sadd a monnin ke california |

Example  "mac store santa monica california"

| LM Scale | Result |
|:---:|:---|
| 0 | mm aq ck storr sade a mon in ck a california |
| 0.01 | maxtor sadd a monnin ke california |
| 1 | maxtor sad a monica california |
| 2 | maxtor sata monica california |

Example  "mac store santa monica california"

| LM Scale | Result |
|---|---|
| 0 | mm aq ck storr sade a mon in ck a california |
| 0.01 | maxtor sadd a monnin ke california |
| 1 | maxtor sad a monica california |
| 2 | maxtor sata monica california |
| 3-19 | mac store santa monica california |

Example   "mac store santa monica california"

| LM Scale | Result |
|----------|--------|
| 0 | mm aq ck storr sade a mon in ck a california |
| 0.01 | maxtor sadd a monnin ke california |
| 1 | maxtor sad a monica california |
| 2 | maxtor sata monica california |
| 3-19 | mac store santa monica california |
| 20 | maxtor pharmacies |

Example   "mac store santa monica california"

| LM Scale | Result |
|----------|--------|
| 0 | mm aq ck storr sade a mon in ck a california |
| 0.01 | maxtor sadd a monnin ke california |
| 1 | maxtor sad a monica california |
| 2 | maxtor sata monica california |
| 3-19 | mac store santa monica california |
| 20 | maxtor pharmacies |
| 25+ | restaurants |

# Summary

Brief tour of the toolbox

Challenges for building ASR systems

A few pointers for getting started

# But what about design?!

That's next!

# Speech technologies for interactive mobile applications – a primer

# Thanks!

Jason D. Williams

at&t

AVIxD / IxD workshop – August 2011