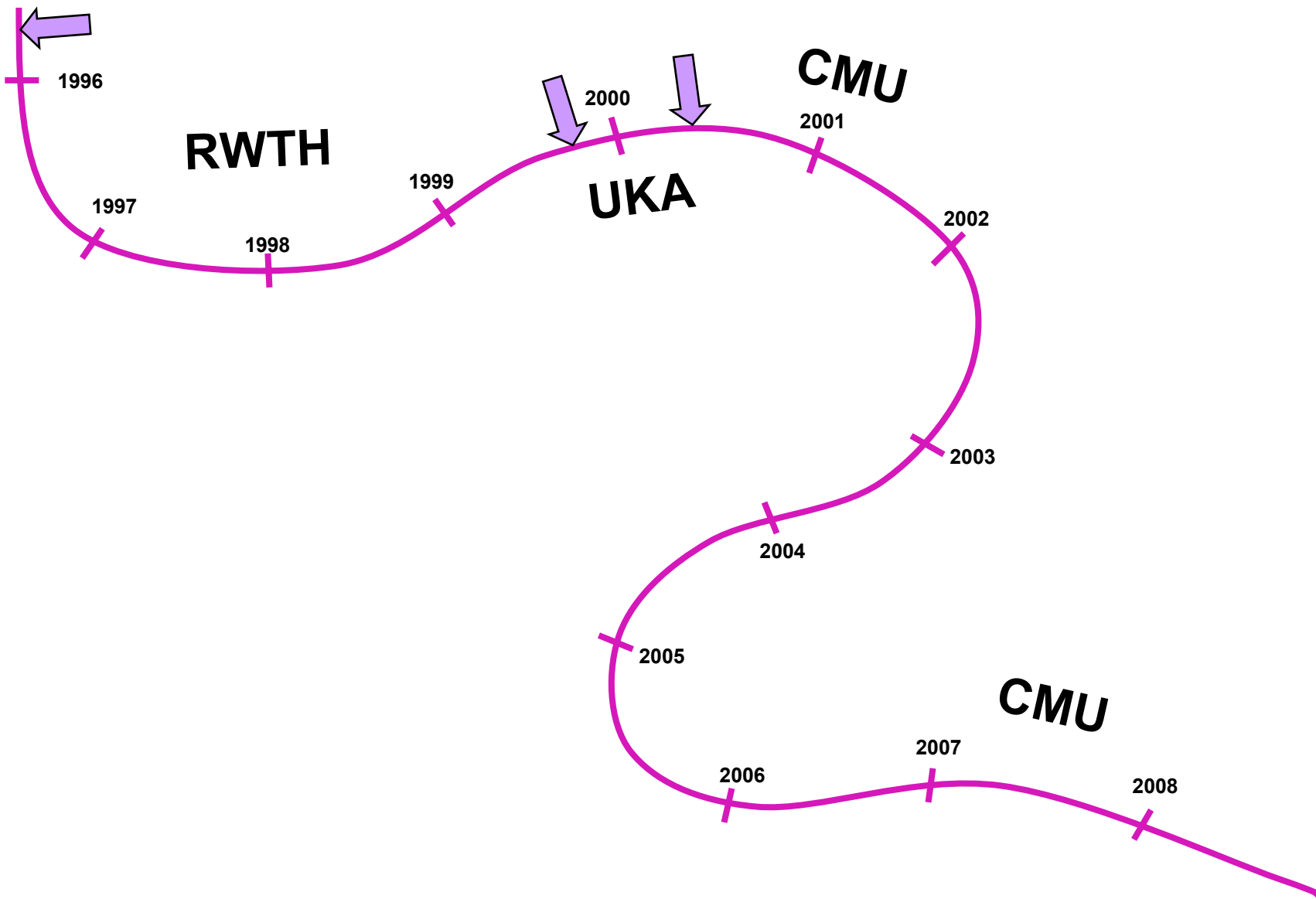# Speech Translation
# From Domain-Limited to Domain Unlimited Translation Tasks

Stephan Vogel

InterACT
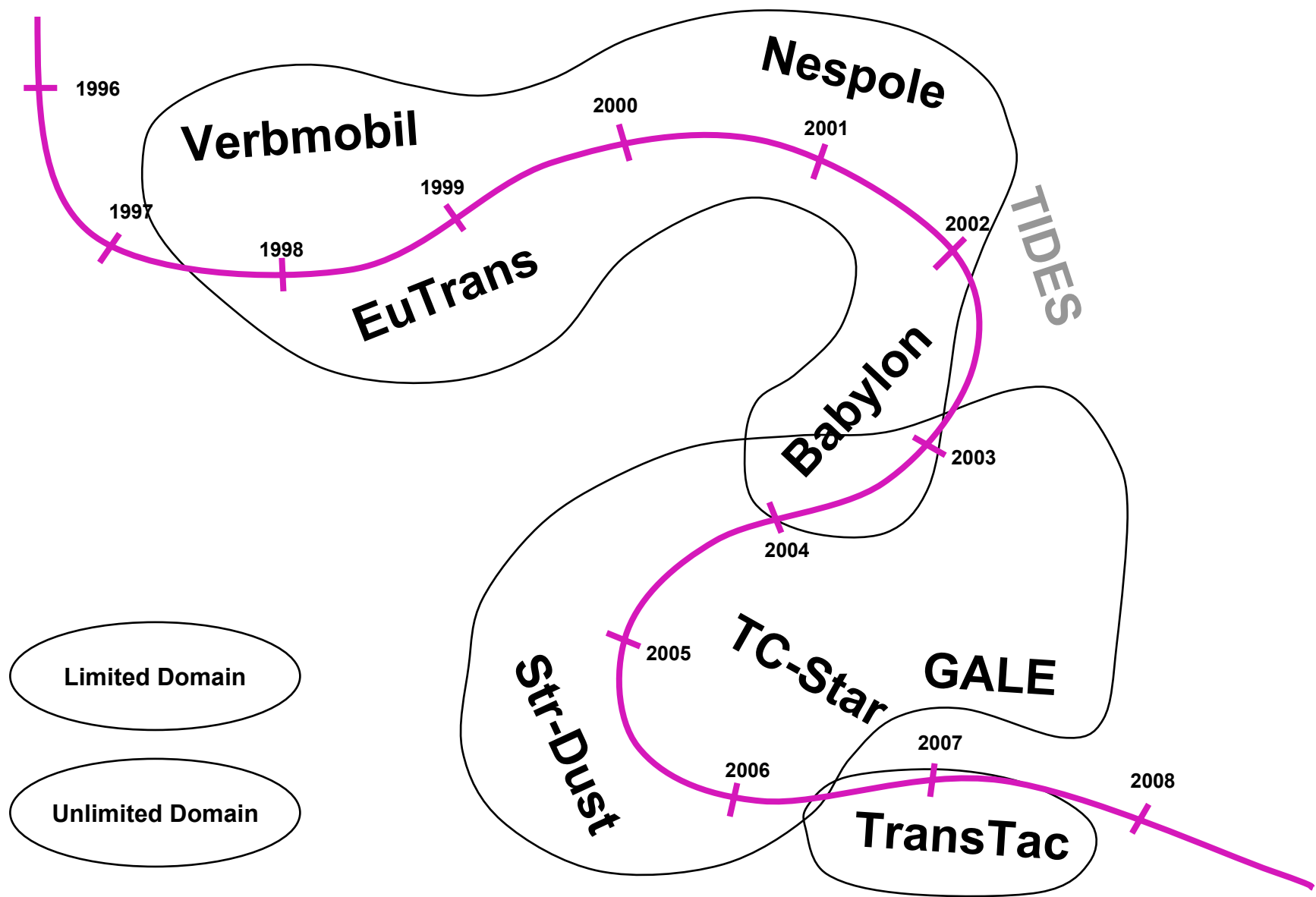Language Technologies Institute
Carnegie Mellon

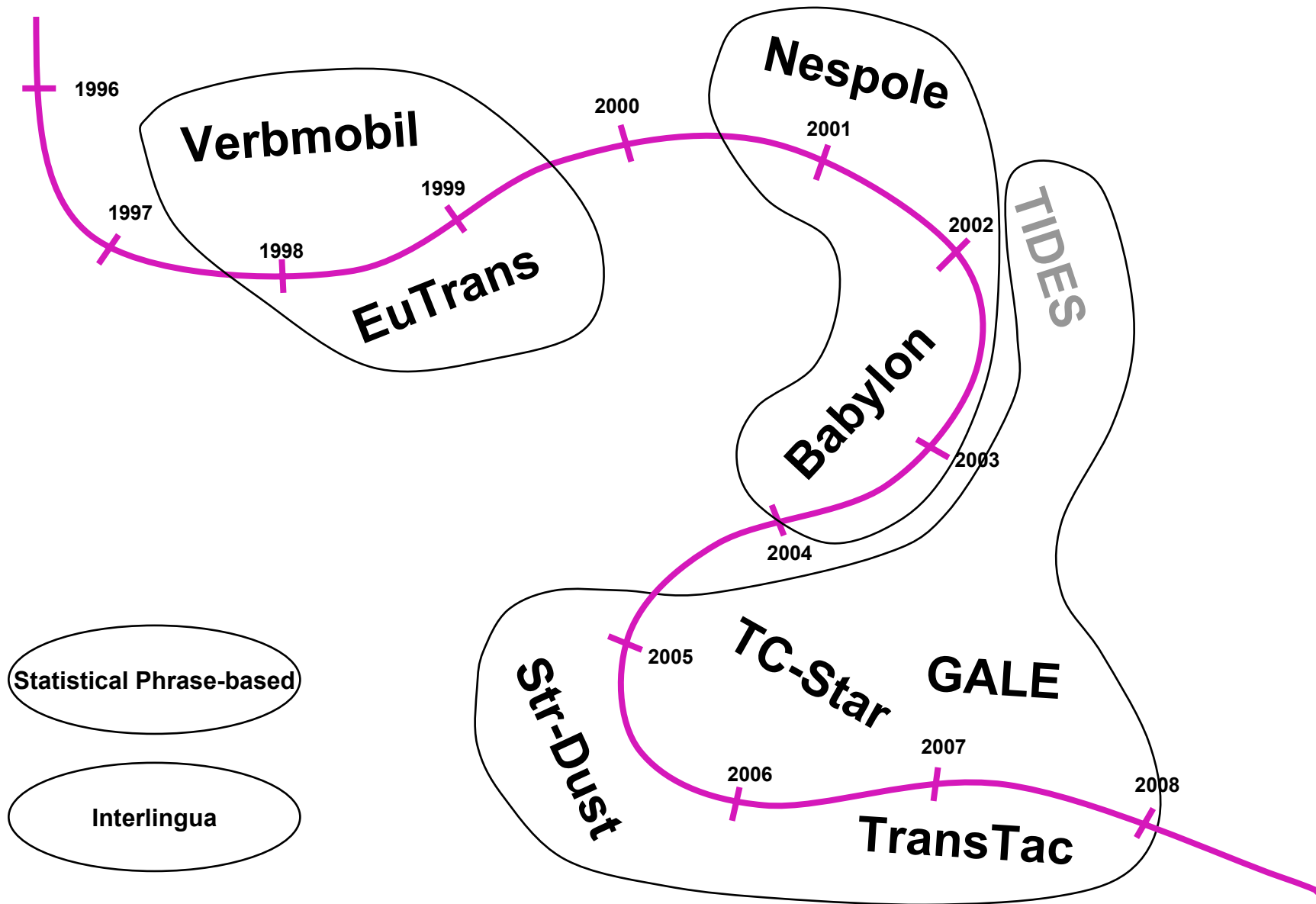# A (Personal) Speech Translation Journey



**1996**

**RWTH**

**2000**

**CMU**

**2001**

**1997**

**1999**

**UKA**

**1998**

**2002**

**2003**

**2004**

**2005**

**CMU**

**2007**

**2006**

**2008**

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# A (Personal) Speech Translation Journey



Verbmobil

Nespole

EuTrans

TIDES

Babylon

Str-Dust

TC-Star

GALE

TransTac

Limited Domain

Unlimited Domain

1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# The Verbmobil Project



Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Verbmobil Highlights

o Multiple recognizers (languages and systems)

o Multiple translation engines
  - o Transfer, with deep syntactic/semantic analysis
  - o Dialog-Act-based
  - o Example based
  - o Statistical

o System combination for translation engines
  - o Did not work

o Prosodic annotation (probabilistic, on lattices)
  - o Segment boundaries on various levels
  - o Question
  - o Accent
  - o Actually used by translation modules

o Disfluency detection and removal (on lattices)

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Nespole!

o Project

  o C-Star partners

  o Interlingua-based translation (English, German, Italian)

  o Travel domain

o Outside of project: test, if SMT is possible on very small corpora

Train: 3182 parallel speech dialog units (sentences)

| Language | English | German |
|---|---|---|
| Tokens | 15572 | 14992 |
| Vocabulary | 1032 | 1338 |
| Singletons | 404 | 620 |

Test: 70 Parallel SDUs

| | German | Reference A | Reference B |
|---|---|---|---|
| Tokens | 437 | 610 | 607 |
| Vocabulary | 183 (45 oov) | 165 | 160 |

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Evaluation

o **Human Scoring**
  - o Good, Okay, Bad (c.f. Nespole evaluation)
  - o Collapsed into a „human score" on [0,1] (good = 1.0, okay = 0.5)
o **Bleu Score**
  - o Average of N-gram precisions from (1..N), typically N=3 or 4
  - o Penalty for short translations to substitute for recall measure
  - o Numbers ranging 0.0 ... 1.0, higher is better

|        |     | Good | Okay | Bad | Score | Bleu  |
|--------|-----|------|------|-----|-------|-------|
| Text   | IF  | 77   | 104  | 227 | 0,32  | 0,068 |
|        | SMT | 127  | 80   | 205 | 0,40  | 0,333 |
| Speech | IF  | 64   | 101  | 243 | 0,28  | 0,059 |
|        | SMT | 95   | 83   | 227 | 0,34  | 0,262 |

o **Results**
  - o SMT works (as good) with very small training data
  - o Both systems pretty useless (OOV!!!)

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Translation System for Travel Domain

Central question:

How to build a translation system which can be used when you are "On the Run"

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Translation System for Travel Domain

o  (Phrase-based) SMT is soooo greedy

o  Memory

  o  Big phrase tables
  o  Large N-gram language model

o  Computation

  o  Decoder needs to generate and
    evaluate many alternative translations

o  Typical machine for SMT systems

  o  2GHz CPU, 4++GB memory
    Linux machine

**Central question:**
  How to build a translation system which can be used
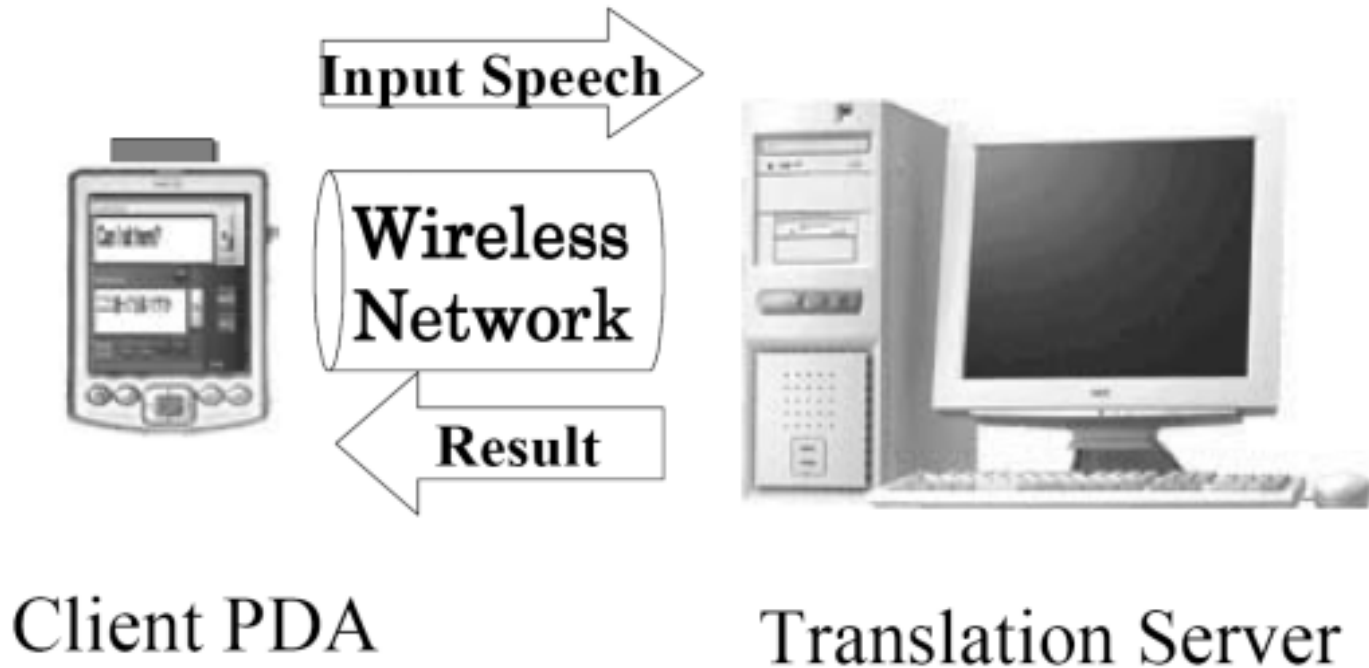  when you are "On the Run"

# Option I

o Travel with powerful computers

# Option II

o  Communicate with a powerful server via wireless*



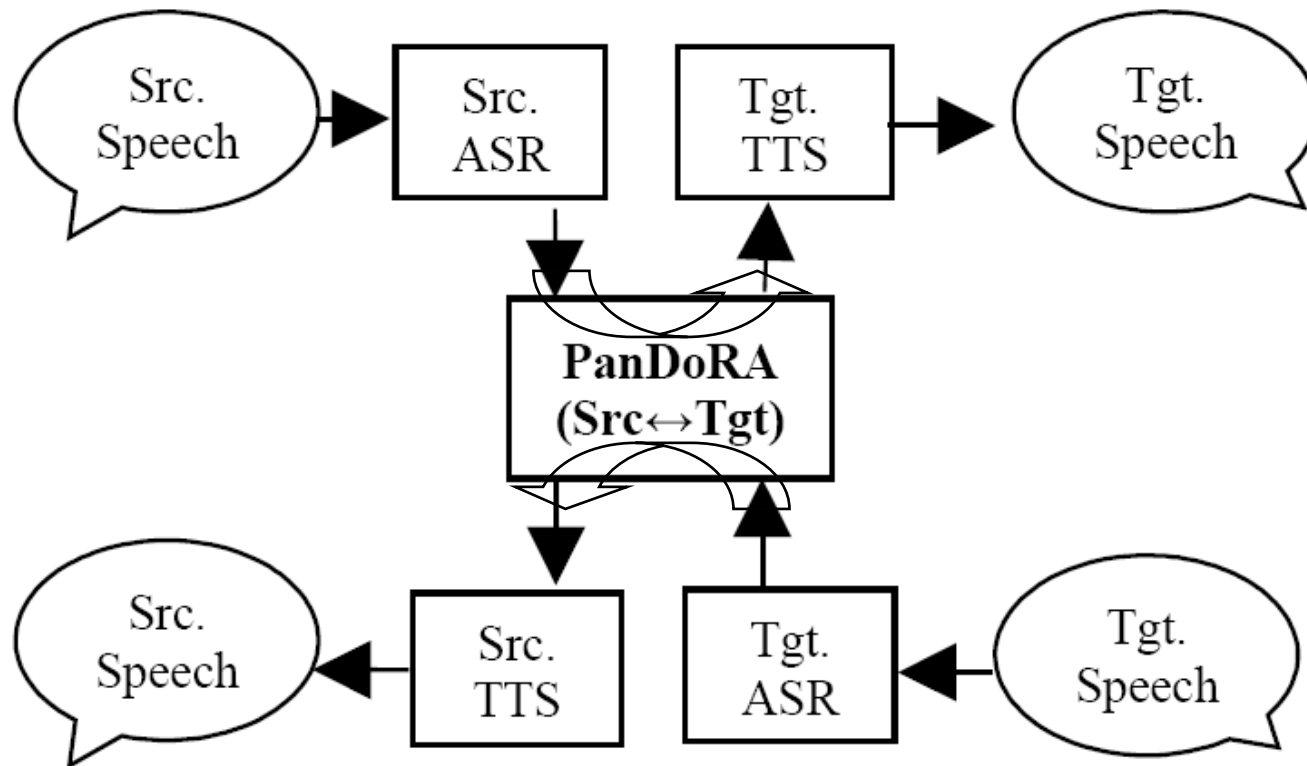o  Problem:  Wireless networks are not always available

---

\* Yamabana et al. 2003. *A speech translation system with mobile wireless clients*.

# Option III: PanDoRA

o  Smart engineering of the SMT system
   so that it can run on hand-held devices
o  All components on device
   o  No dependency on connectivity

# PanDoRA: System Architecture

# Challenges

o Memory

   o SMT systems require several GB RAM on PCs

   o Hand-held devices have very small dynamic memory

   o 64MB available on iPaq 2750; shared with ASR and TTS

o CPU

   o Speech translator requires real-time translations

   o CPUs on hand-held devices are usually slow (e.g. 600M Hz) and

   o There are no coprocessors for floating-point arithmetic

# Solutions

o   Compact data structures

o   Integerized computation

o   Efficient decoding
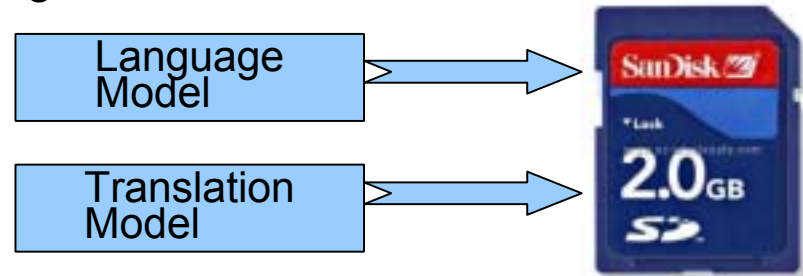
o   Minimum on-device computation

# Compact Data Structure

o Each phrase is represented by its <location,length> in the corpus rather than a string of words

o A translation pair is represented by <src phrase index, tgt phrase index, score>

احتفال المدرسة | # school festival is # 0.0034
احتفال المدرسة | # school festival is hold # 0.0031
احتفال المدرسة عقد # school festival is hold # 0.9980
احتفال دل | # dolls' festival # 0.5431
احتفال دل للفتيات | # dolls' festival # 0.3999
احتفاليا| # festive # 0.7535

<100, 200, 0.0034>
<100, 201, 0.0031>
<101, 201, 0.9980>
<102, 202, 0.5431>
<103, 202, 0.3999>
<104, 203, 0.7535>

**Arabic corpus:**
… احتفال المدرسة
احتفال المدرسة عقد ..

**English corpus:**
this year's school festival is hold on September .
….
…. there will be a doll's festival….
…festive….

# Compact Data Structure

o Data records now have fixed size in bytes, independent of the length of a phrase.

o Sorted list of all model data is stored on external memory card
  o Much larger capacity than SDRAM. e.g. 2GB SD card
  o Saves SDRAM for decoding

o Search for matched phrases by binary search in the sorted list

o Access the information by 'seek'ing the record from a file by record id.
  o e.g. getPhrase(id); getTranslationOfPhrase(id)...

o With compact data structure a 5.6 million entry phrase table requires 65MB on disk

o Specification for current implementation
  o 64K Vocabulary for each language
  o Up to 256 million unique src/tgt phrases allowed
  o Up to 4 billion phrase pairs allowed

Language Model

Translation Model

# Integerized Computing

o SMT systems are "translating by numbers"

o Hand-held devices usually do not have numerical co processors

o Slow with floating points calculations

o PanDoRA's solution: integerized computing

    o Convert all probabilities to cost = -log(prob)

    o Quantize the cost to integer bins between [0, 4095]

    o Probabilities put into the same bin are considered equal during decoding.

| TM | LM | BLEU | NIST |
|---|---|---|---|
| Float | Float | 19.87 | 8.03 |
| Int | Float | 19.82 | 7.99 |
| Float | Int | 19.94 | 8.03 |
| Int | Int | 19.93 | 8.04 |

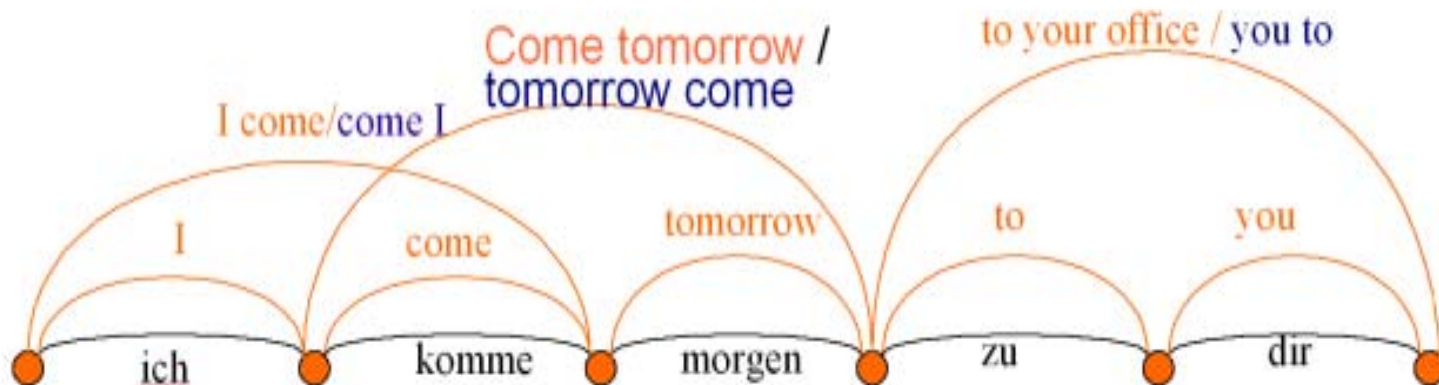# Minimum On-device Computation

o Log-linear model in SMT
  o Multiple model features

$$P(e_1^I \mid f_1^J) = \frac{\exp(\sum_{m=1}^{M} \lambda_m \phi_m (e_1^I, f_1^J))}{Z}$$

o Off-line model optimization
  o Minimum Error (MER) training off-line
  o Combine weighted translation model scores to a single score

# Decoding

- o Inverted Transduction Grammar (ITG) style decoding [bottom-up]
  - o Translation as parsing
  - o Inverted Transduction Grammar (ITG) (Wu, 1997)
    - o Straight transduction: X-><f1f2 | e1e2>
    - o Invert transduction: X-><f1f2 | e2e1>
  - o Combine adjacent partial hypotheses either straight or inverted to create new hypothesis covering the combined source range.
  - o Allows for long distance reordering. Effective for language pairs with dramatically different word ordering. E.g. Japanese/English
  - o Language model adjusted on the boundary words.

# Experiments on Japanese/English

|  | Japanese | English |
|---|---|---|
| Word Tokens | 1.2M | 1.0M |
| Word Types | 18K | 13K |
| Sentences | 162K | 162K |
| Avg. Sent. Len. | 7.32 words | 6.18 words |

Statistics of the BTEC Jp/En training data

| J→E Phrase Pairs | 4,648,018 |
|---|---|
| E→J Phrase Pairs | 4,871,862 |
| Uniq. Japanese Phrases | 1,396,719 |
| Uniq. English Phrases | 1,015,821 |

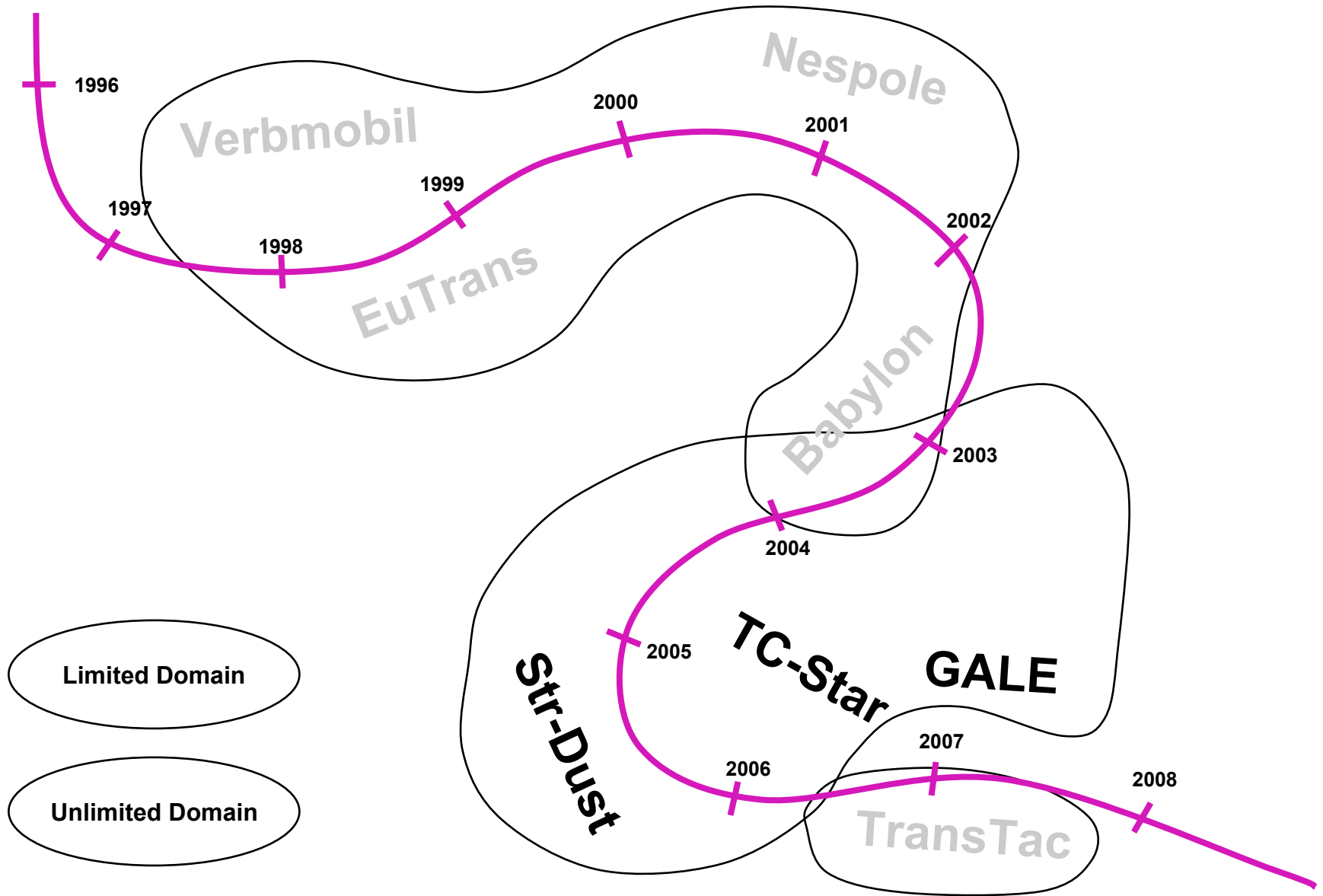Japanese↔English phrase translation model.

| Phrase Table | LM | Reorder | STTK | PanDoRA |
|---|---|---|---|---|
| Pruning | SRI | No | 46.2 | 45.93 |
|  | 3-gram | Yes | 52.4 | 54.59 |
|  | SALM | Yes | 53.6 |  |
| No Pruning | SRI | No | 50.3 | 49.96 |
|  | 3-gram | Yes |  | 58.64 |
|  | SALM | Yes | 59.1 |  |

Translation results of the PanDoRA system on the Japanese to English task, compared with the performance of STTK.

# Travel Domain Translation Demo



Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# From Limited towards Unlimited



Limited Domain

Unlimited Domain

1996

1997

1998

1999

2000

2001

2002

2003

2004

2005

2006

2007

2008

Verbmobil

Nespole

EuTrans

Babylon

Str-Dust

TC-Star

GALE

TransTac

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# TC-Star

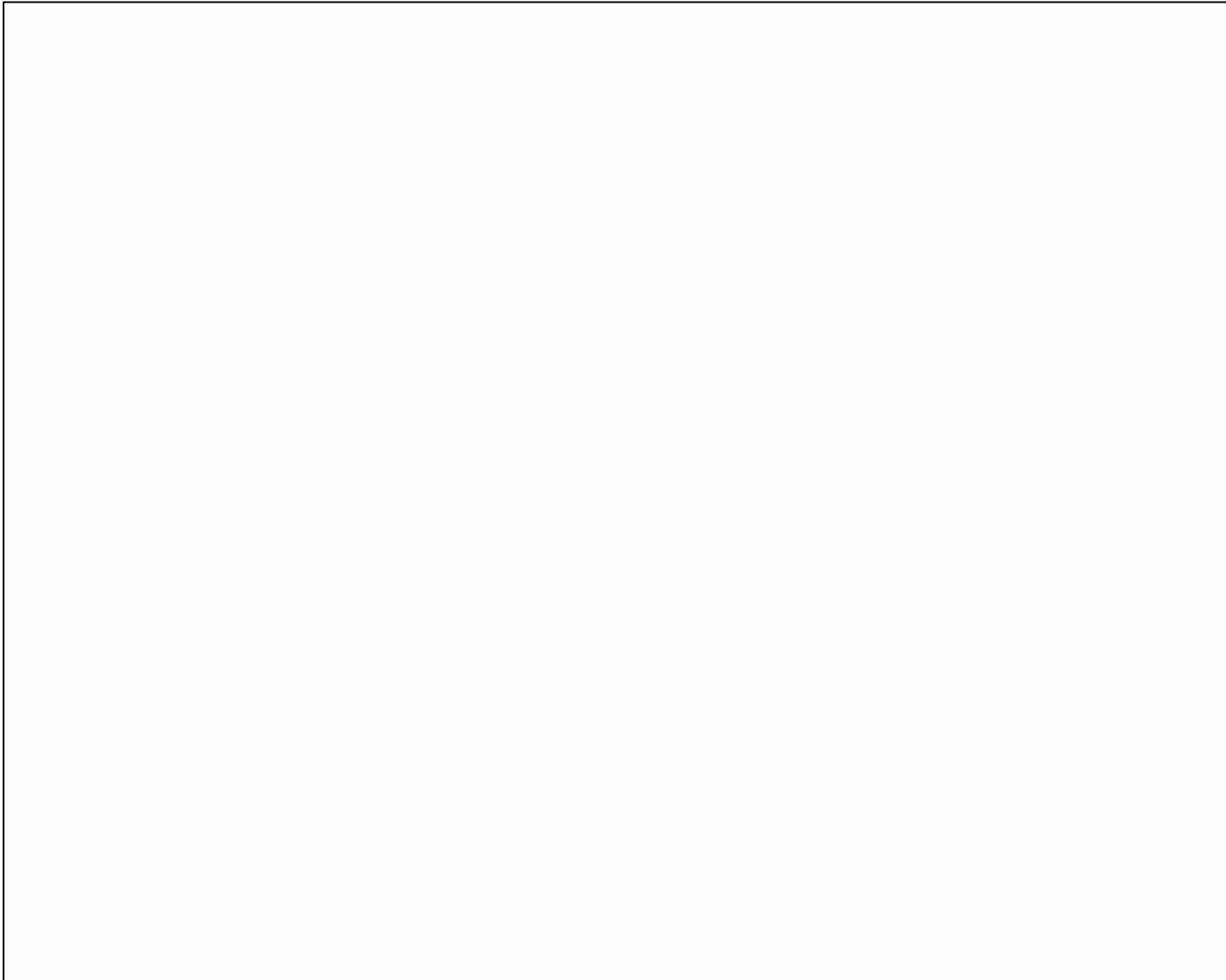**Technology and Corpora for Speech to Speech Translation**

o Financed by European Commission within the Sixth Program

o Research in all core technologies for Speech-to-Speech Translation

o Long-term research goals of the project:

   o Effective Spoken Language Translation of unrestricted conversational speech on large domains of discourse.

   o Speech recognition able to adapt to and perform reliably under varying speaking styles, recording conditions, and for different user communities.

   o Effective integration of speech recognition and translation into a unique statistically sound framework. (Lattice translations as explicit goal)

   o General expressive speech synthesis imitating the human voice. (Development of new models for prosody, emotions and expressive speech.)

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# TC-Star Evaluations

o EPPS:  European Parliamentary Plenary Sessions)

- o EU >20 official languages
- o Speeches given in native language
- o Simultaneous interpretation (verbatim transcript)
- o Final text edition for publication (often significantly edited)

o TC-Star Evaluations

- o Show-case for potential translation (interpretation) service in the European Parliament
- o Spanish-to-English and English-to-Spanish
  - o Verbatim transcript
  - o ASR output (1-best, lattice)
  - o Final text edition
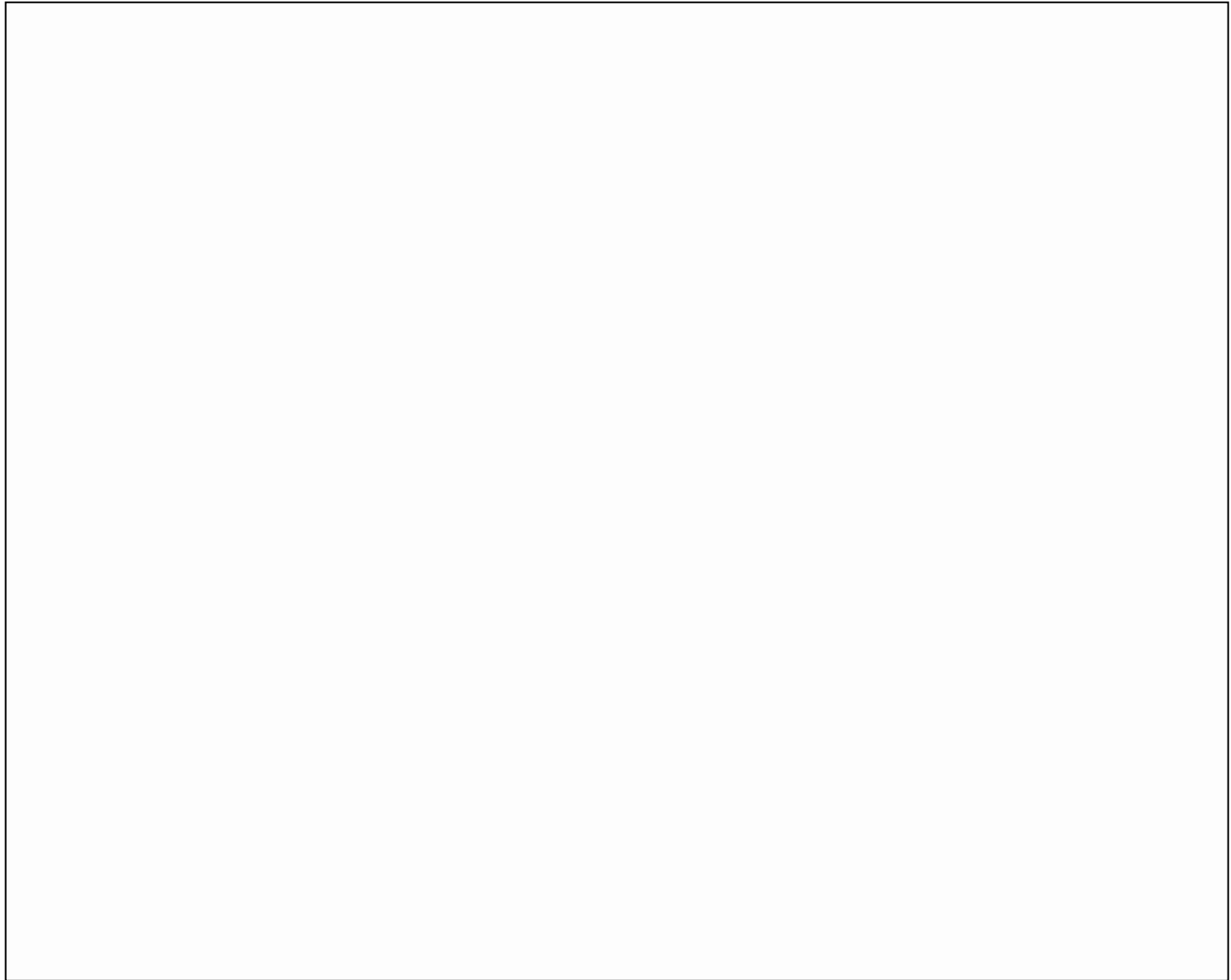- o Also Chinese-to-English (NIST MT eval clone?)

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# TC-Star Demo

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Lecture Translator

o Speech recognition

- o Large vocabulary
- o Continuous speech
- o Fair amount of disfluencies

o Translation

- o English -> Spanish
- o English -> German
  - o Trained on EPPS (European parliament) data
- o English -> Arabic
  - o Trained on UN corpus
- o Additional LM data for adaptation to domain

o Application/Demo

- o Translating lectures, e.g on speech and translation technology

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Lecture Translation Demo

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Lecture Translation Demo

MORNING THANK YOU VERY MUCH IT IS VERY E NICE TO HAVE YOU HERE IT IS A PLEASURE TO
TALK TO YOU ABOUT CERTAIN TECHNOLOGICAL ADVANCES THAT WE WOULD LIKE TO SHOW YOU
ON THIS DAY FOR THE FIRST TIME
WE'D LIKE TO INTRODUCE TO YOU
SPEECH

mañana muchas gracias es muy bonito decir que usted aquí es un placer
le voy a hablar algunas tecnológicos nos gustaría mostrarle
en este día de la primera vez
nos gustaría introducir les

heute vormittag vielen dank es ist sehr e zu haben wenn sie in nizza hier es ist eine freude
sprechen sie über bestimmte technische daß gerne zeigen würden wir fortschritte
an diesem tag zum ersten mal
wir möchten ihnen gerne vorstellen

الصباح اشكركم كثيرا جدا ان ه والتحاور هل هو هنا ويسرني ان
التحدث عن لكم بعض التكنولوجي نود ان نبين لكم
على

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Translation Server for Informedia

o **Informedia system captures videos**

o **Sound track is send to ASR server**

o **1-best result send to translation server**

o **Results archived in Informedia database**

o **Highlighting in text synchronized with video**

..¥Demos¥InterACT-BNT¥InterACT-BNT.exe



Recognized Arabic

Translation into English

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Some Research Topics

o Tight coupling between recognition and translation
   o Lattice translation
   o Using ASR features (LM, acoustic scores, confidences) in MT decoder
   o End-to-End optimization

o Segmentation and restoring punctuation

o Disfluency detection and removal

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Lattice Translation: Motivation

o Lattice translation is used for tighter coupling between speech recognition and translation

- o Use alternative word sequences
- o Lattice is more efficient than n-best list
- o Word lattice has lower word error rate than 1-best hypothesis -> perhaps we find a better path in translation
- o Perhaps a different path gives better translation even when more recognition errors
- o Decouple processing steps, yet use relevant information for end-to-end optimization

o Lattices can be used to encode alternatives arising from other knowledge sources

- o Disfluency annotation: add edges to lattice to skip over disfluencies
- o Add synonyms and paraphrases
- o Add morphological variants, esp. for unknown words
- o Allow different word segmentation (e.g. For Chinese)
- o Partial word reordering

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Requirement: Translating Lattices

o  Connection between speech recognizer and translation system can be

   o  1-best hypothesis

   o  N-best hypotheses list

   o  Word lattice

   o  Confusion lattice

o  Why lattice

   o  Compact representation of many hypotheses

   o  Word lattice has lower word error rate than 1-best hypothesis
      -> perhaps we find a better path in translation

   o  Perhaps a different path gives better translation even if it has more word errors

   o  Can use results from consolidation modules without making a hard decision: consolidation module adds paths to the lattice

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Requirement: Using ASR Features

o Problem when translating lattices
- o Perhaps bad path is chosen because it is "easier" to translate
- o Shorter paths are more likely to be selected

o Solution: Use Recognizer Features
- o Acoustic scores: local to source word
- o Language model scores: non-local
  - o Expansion of the word lattice for different histories
  - o Calculation on the fly, i.e. in the SMT decoder
  - o Affected by word reordering
- o Confidence scores

o Other features might be interesting
- o Segmentation: probabilities for boundary after each word
- o Consolidation modules: probabilities for new edges

o Goal: overall optimization of end-to-end system

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Requirement: Word Reordering

o Long distance reordering
  - o Allow first word to be translated last, last word to be translated first
  - o All permutations is too expensive
  - o Typically only a small number of reorderings in a sentence
o Types of reordering
  - o Jump ahead a few words, fill gap one or two steps later (only local word reordering)
  - o Leave small number of gaps, fill at any time (global reordering possible)
  - o Binary Tree with swapping sub-trees: not all reorderings orderings possible, but pretty much all which appear in natural languages (global reordering possible)

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007
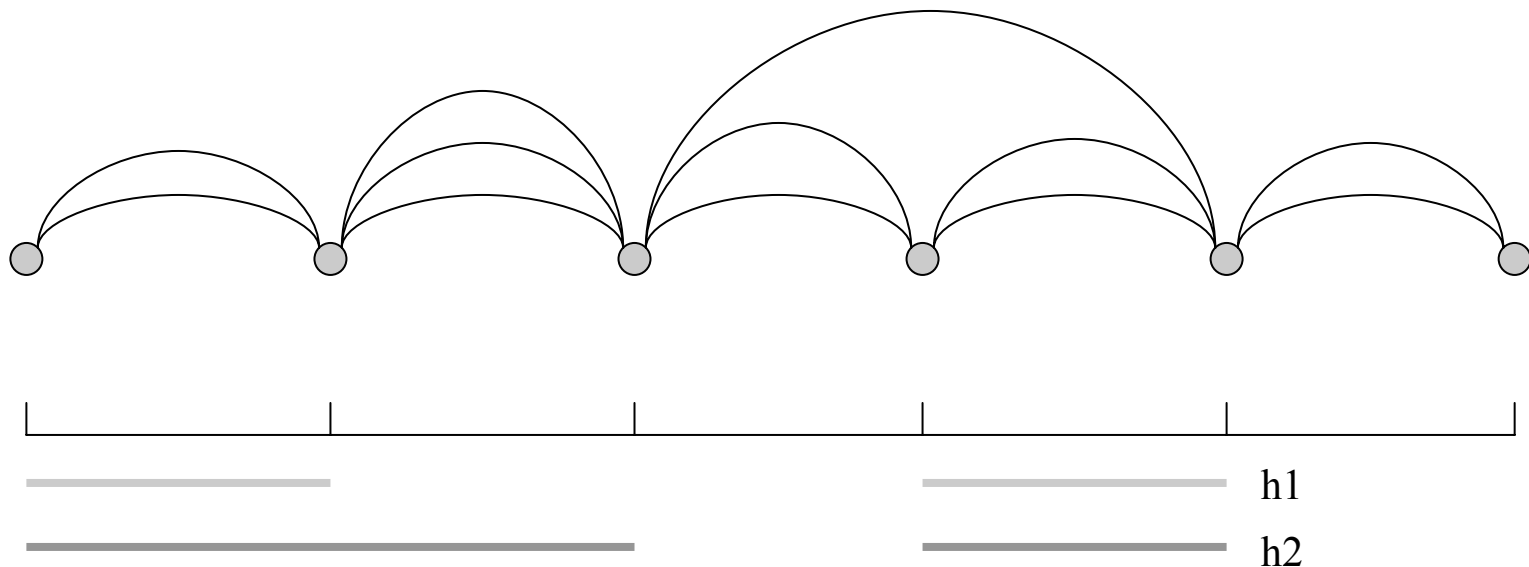
# Different Reordering Strategies

o **All permutations**
  - o Any re-ordering possible
  - o Complexity of traveling salesman -> only possible for very short sentences

o **Small jumps ahead – filling in the gaps pretty soon**
  - o Only local word reordering

o **Leaving small number of gaps – fill in at any time**
  - o Allows for global but limited reordering in one direction
  - o Similar decoding complexity – exponential in number of gaps
  - o IBM-style reordering (described in IBM patent)

o **Merging neighboring regions with swap – no gaps at all**
  - o Allows for global reordering in both directions
  - o Complexity lower than 1, but higher than 2 and 3

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Dealing with Word Order – Coverage Info

o Need to know which source words have already been translated
  - o Don't want to miss some words
  - o Don't want to translate words twice
  - o Can compare hypotheses which cover the same words

o Coverage vector to store this information
  - o For 'small jumps ahead': position of first gap plus short bit vector
  - o For 'small number of gaps': array of positions of uncovered words
  - o For 'merging neighboring regions': left and right position

o All simple and easy for sentences as input – but does this work for lattices?

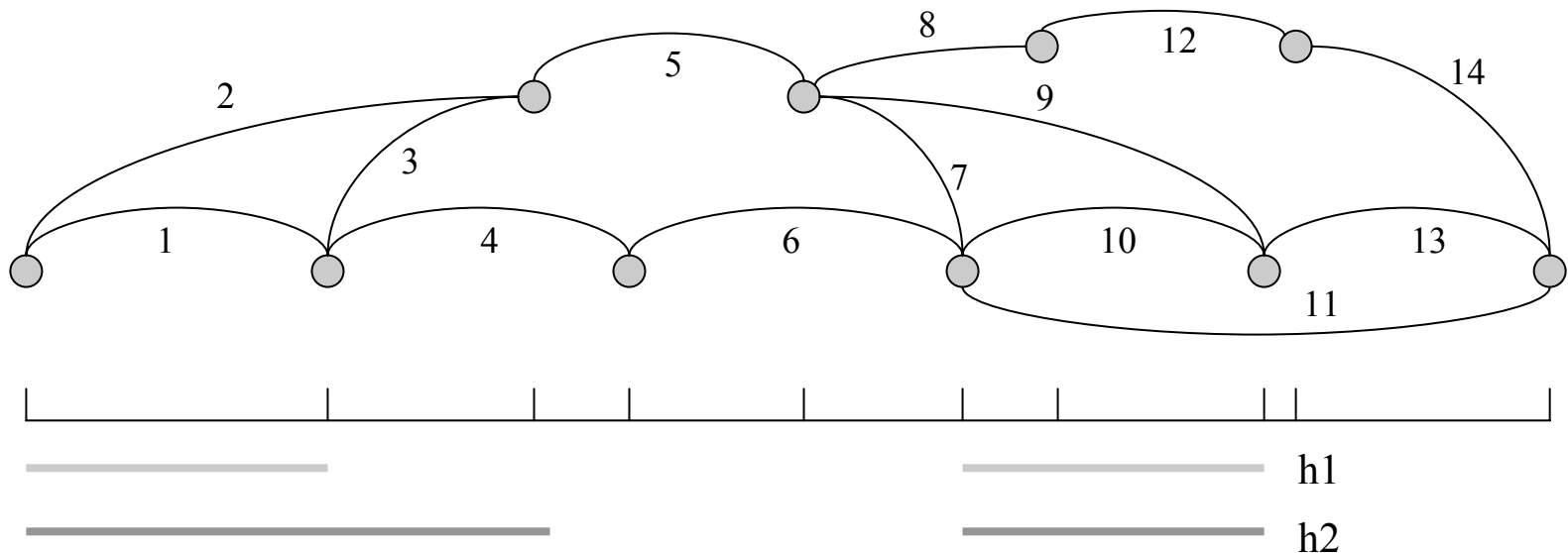Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Coverage Information – Confusion Nets

o Empty edges can be eliminated by expanding all incoming edges to next node

o Structure of confusion net is same as structure of translation lattice after inserting word and phrase translations

o Works fine for all reordering strategies

h1

h2

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Coverage Information – General Lattices

o Reordering strategy 1 is somewhat messy
  - o Can be simulated with strategy 2, so don't bother
o Strategies 2 and 3 pose no additional problem
  - o For 2: Recombine hyps when ending in same node and having same untranslated edges
  - o For 3: Recombine hyps which have same left and right node

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Using Source Language Model

o   Edges for target words/phrases need to be linked to source word edges

o   Easiest in reordering strategy 3
   o   Run over all words attached to source edges, perhaps
   o   Run over all words in target edges
   o   One side could have the two joint segments swapped
o   In reordering strategy 2
   o   Calculate source LM prob when leaving word/phrase untranslated
   o   Calculate target LM prob when filling the gap

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Lattice Translation Results

o A number of papers
  - o General word graphs
  - o Confusion nets
  - o N-best lists (with N=10…100) converted into lattices
o Improvements (small and smaller) reported

o General Picture?
  … there is none

  - o Worked last time year - not this year (Wade Shen IWSLT 05 and 06)
  - o Works only on small lattices - works on large lattices (Saleem 04, Matusov 07)
  - o Works only for mid-range WER – works best for low WER (but your WER ranges might be different from mine)
  - o Nice improvements - only small improvements
  - o Translation system selects path with lower WER – selects paths with higher WER (Lane, Paulik, inhouse results)

o Problem:  our models are too weak to select the few better paths in a word lattice among the zillions of worse paths

# Segmentation

o Wrong segmentation introduces errors
  - o Reordering over sentence boundaries
  - o Loosing LM context
  - o Loosing phrase matches

o Segmentation approaches
  - o Simply use pause longer then fixed threshold
  - o Add LM
  - o Add prosodic features

o General picture
  - o Wrong segmentation hurts (our TC-Star eval disaster)
  - o Longer or shorter segments better?
    - o Conflicting results
    - o Longer seems better
  - o Looking at phrase table seems to help
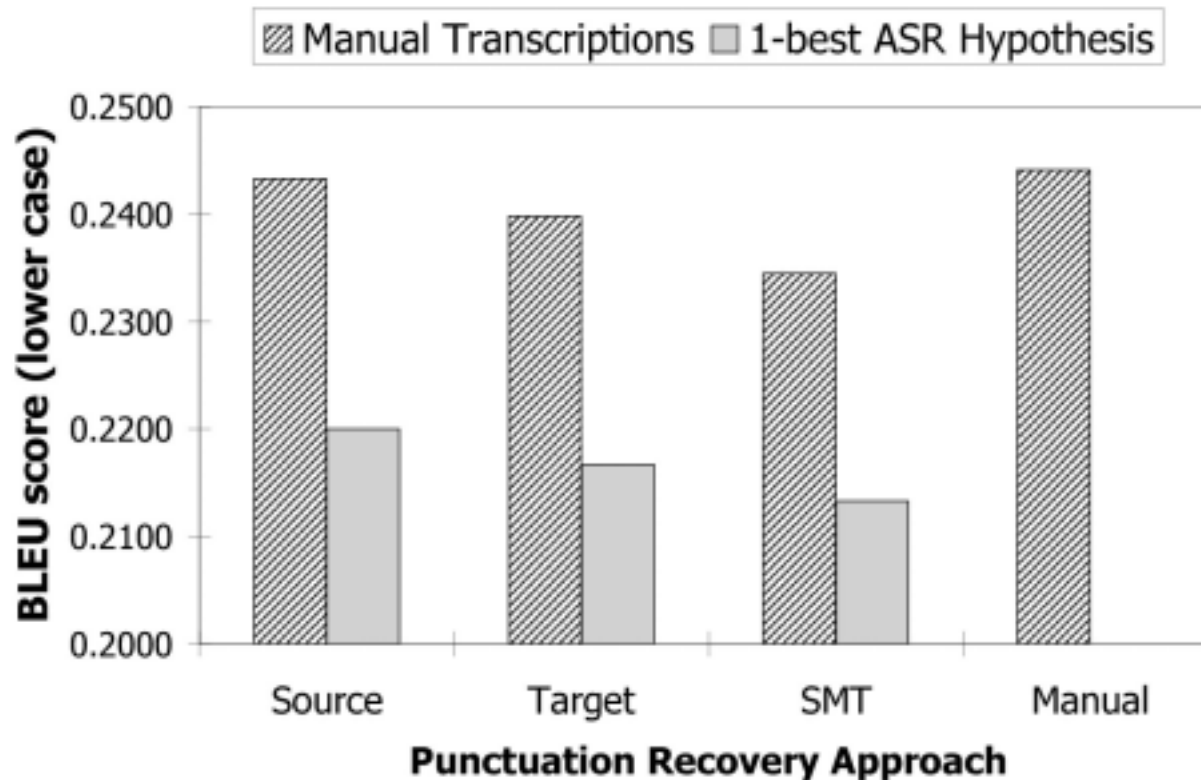
# Punctuation Recovery

Different approaches

o   Use prosodic features as punctuation substitute

o   Hidden event LM
  o   Delete all punctuation from MT training corpus, restore punctuation on target side through HELM
  o   Restore punctuation on source side and used MT trained on corpus with punctuation

o   Asymmetric translation model
  o   Remove punctuation on source side of training data
  o   Or from source side of phrase table
  o   Generate punctuation on target side

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Punctuation Recovery: Some Results

o IWSLT 2007, Japanese-English translation system
o Recover on source side with HELM

|  | Precision | Recall | F-score |
|---|---|---|---|
| **Manual transcripts** | | | |
| Sentence Boundary | 97.8% | 96.8% | 97.3% |
| Secondary Punctuation | 82.1% | 44.2% | 57.5% |
| **1-best ASR Hypothesis:** *Character Error Rate=10.4%* | | | |
| Sentence Boundary | 96.4% | 95.9% | 96.2% |
| Secondary Punctuation | 71.8% | 43.6% | 54.3% |
| **Lower Bound:** *assume sentence boundary at end of utterance* | | | |
| Sentence Boundary | 100% | 63.9% | 77.9% |

# Punctuation Recovery: Some Results



**Source:** Punctuation recovery applied to input source
**Target:** Punctuation recovery applied after translation
(secondary-punctuation ignored during translation)
**SMT:** Punctuation recovered during translation
**Manual:** Input source manually annotated with punctuation

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# User Problems

o Output of speech translation systems is difficult to understand
  o Translation quality
  o Missing or incorrect segmentation: Really big problem for user

o Text display
  o Difficult to 'parse' the system output
  o Text display without punctuation is difficult to read
  o Wrong line-breaks have to be (mentally) undone

o Speech synthesis
  o Rather monotone and unstructured
  o Overlay with original speech
  o Real-time processing by user is hardly possible

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# Next Grand SLT Challenge: Multilingual Meetings

o Several meeting projects: AMI and AMIDA, CHIL, ICSI, …

o So far not translation involved, but would be a natural extension

- o Participants have often different native languages
- o There could be a real need, even more than for lecture translation

o Interesting research issues

- o Latency in translation
  - o How to achieve short latency
  - o How does is affect translation quality
  - o What is 'best' latency for user
- o How to integrate external knowledge sources
  - o Slides, notes from previous meetings
  - o Online search for model adaptation
- o Learning from human interpreters
  - o Strategies to translate despite of missing information
- o Presenting the translations
  - o Audio or video display
  - o How to deal with cross talk and parallel discussions
  - o Tele- and videoconferences seem to provide an easier setup

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007

# My Personal (Speech) Translation Journey ???



**RWTH**

**UKA**

**CMU**

**CMU**

1996
1996
1997
1997
1998
1998
1999
1999
2000
2000
2001
2001
2002
2002
2003
2003
2004
2004
2005
2005
2006
2006
2007
2007
2008
2008
2010
2020
2030
2040
2050

Vogel: Speech Translation, ASRU, Kyoto, Dec 2007