

# NPFL099 - Statistical dialogue systems

## User simulation

Filip Jurčiček

Institute of Formal and Applied Linguistics  
Charles University in Prague  
Czech Republic

Home page: <http://ufal.mff.cuni.cz/~jurcicek>

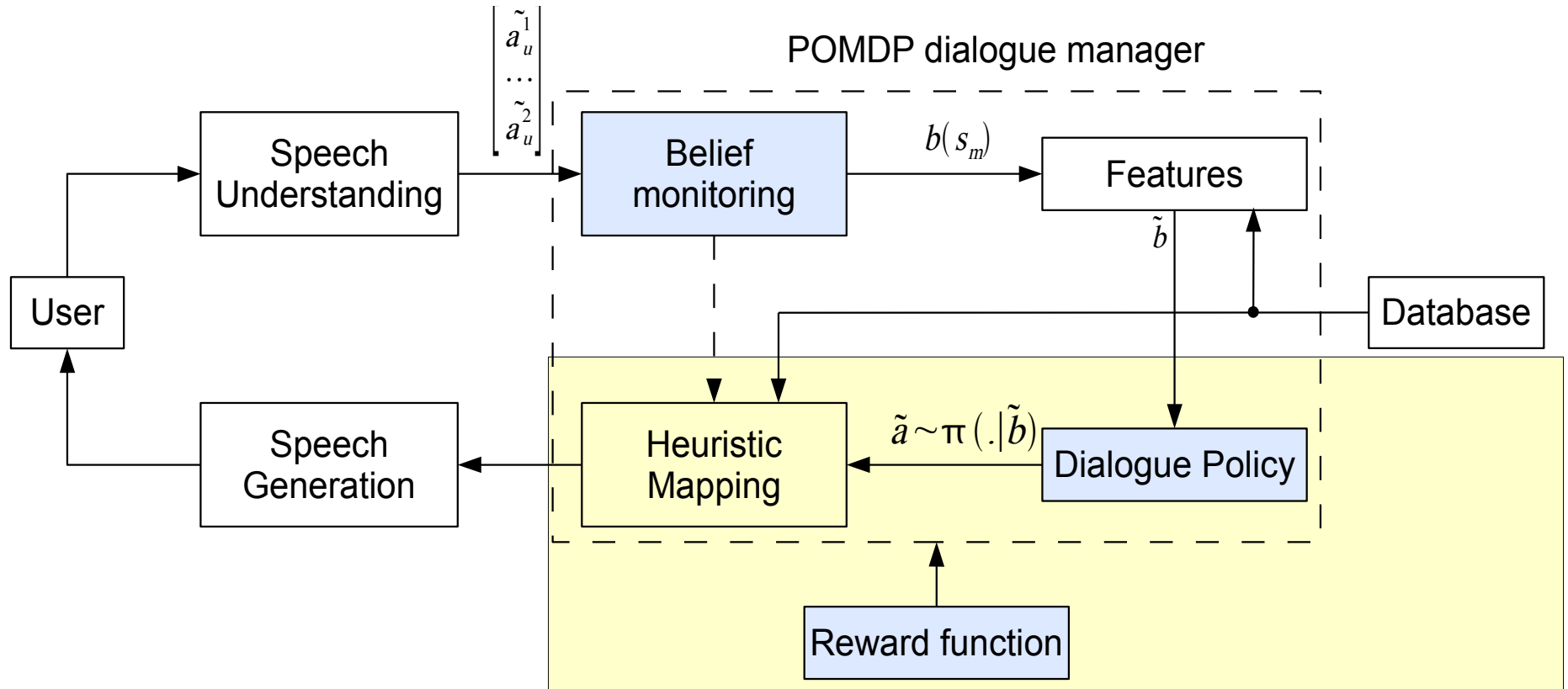
Version: 14/05/2013



# Outline

- User simulation
- Dialogue act level simulations
  - N-gram dialogue act model
  - Agenda based simulation
  - ISU based approach
  - Bayesian user simulation
- Word level simulation
- Speech level simulation

# POMDP Dialogue systems



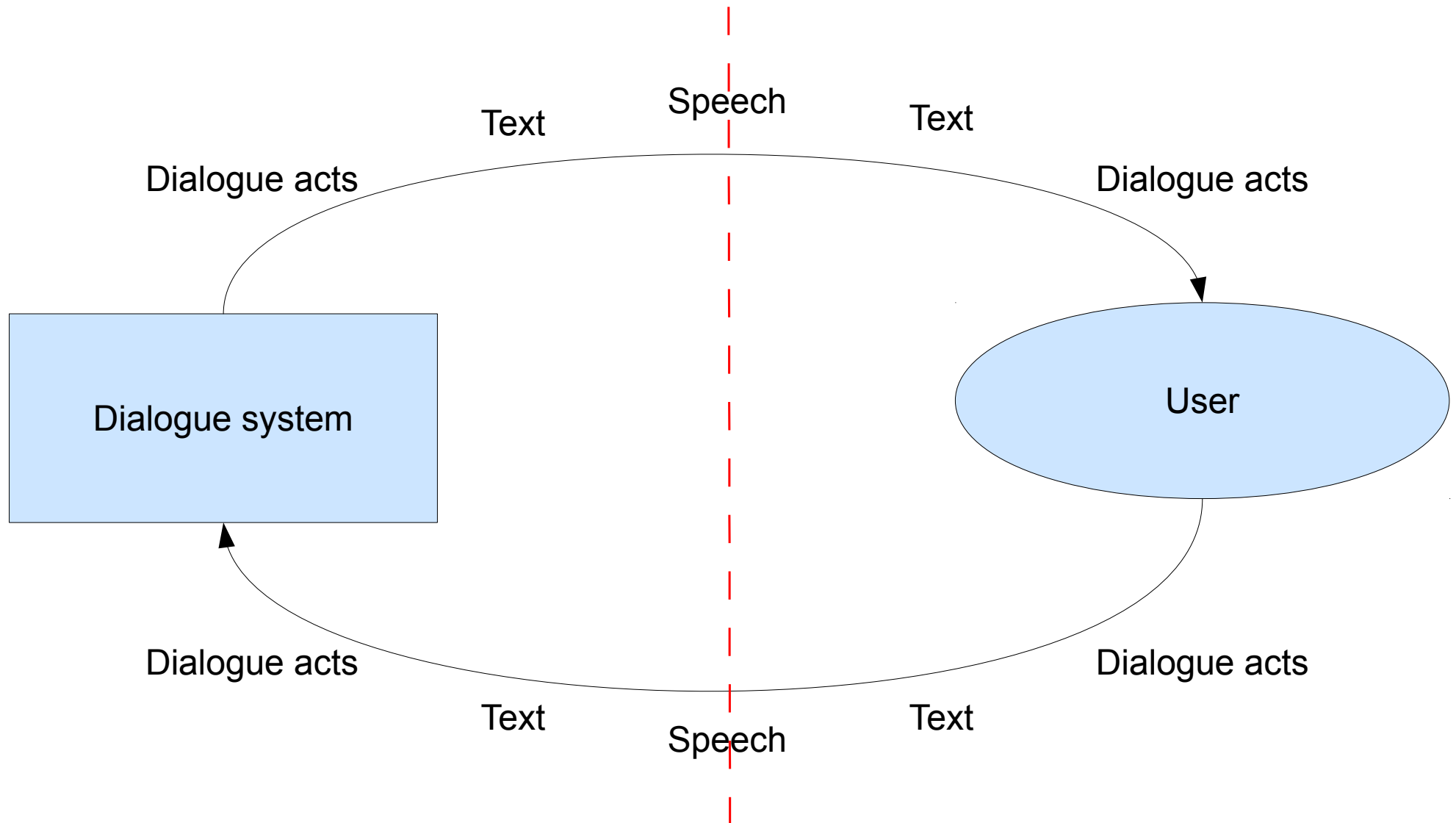
- Ideally the policy is optimized to maximize the reward function

# User simulation

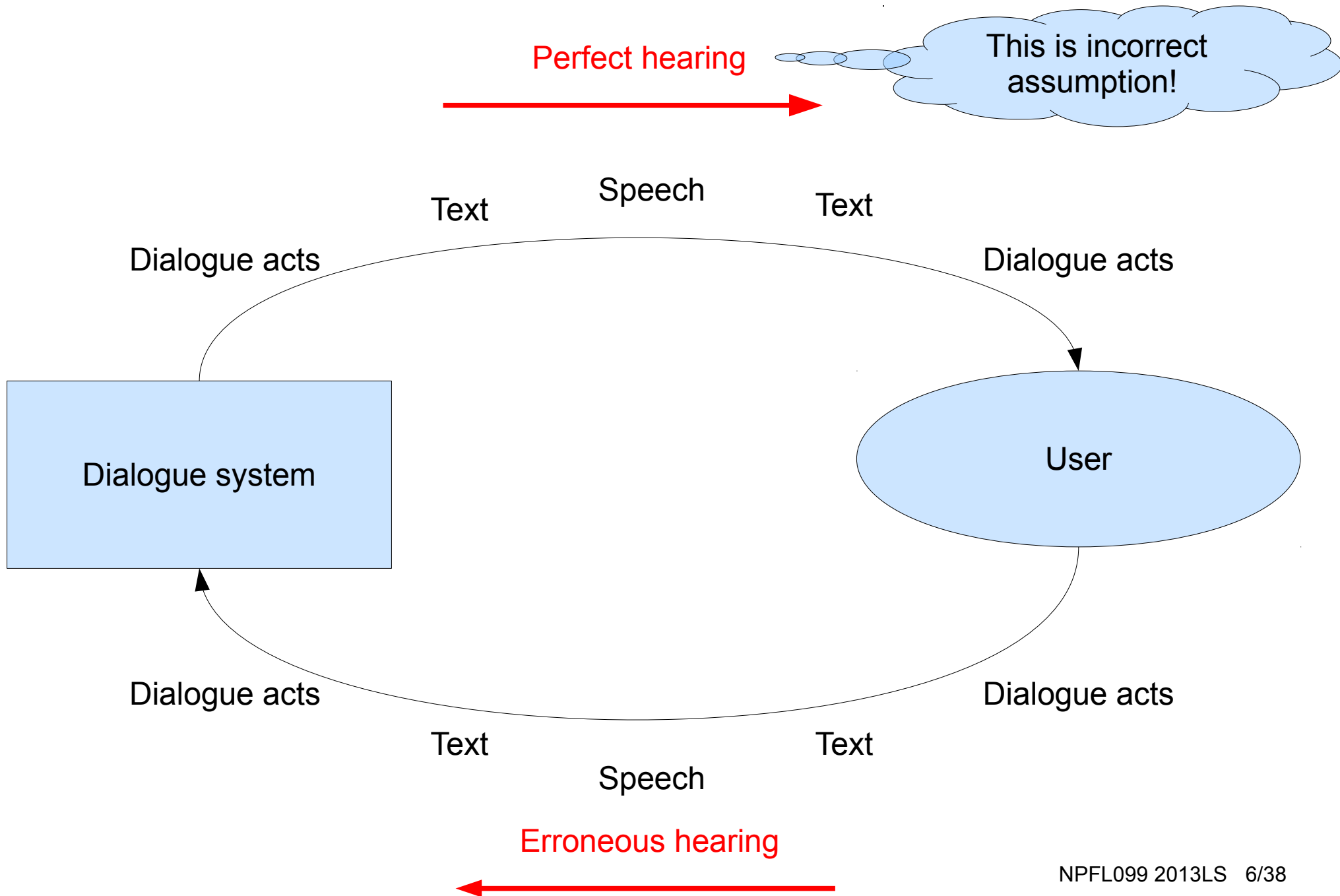
- Ideally, the POMDP dialogue systems would be optimised in interaction with real users
- The problem
  - state-of-the-art techniques still needs to more than 10000 dialogues
- User simulators are used to train and test POMDP techniques
- User simulators are mostly hand-crafted, though parametrised

# User simulation

- It is like building another SDS

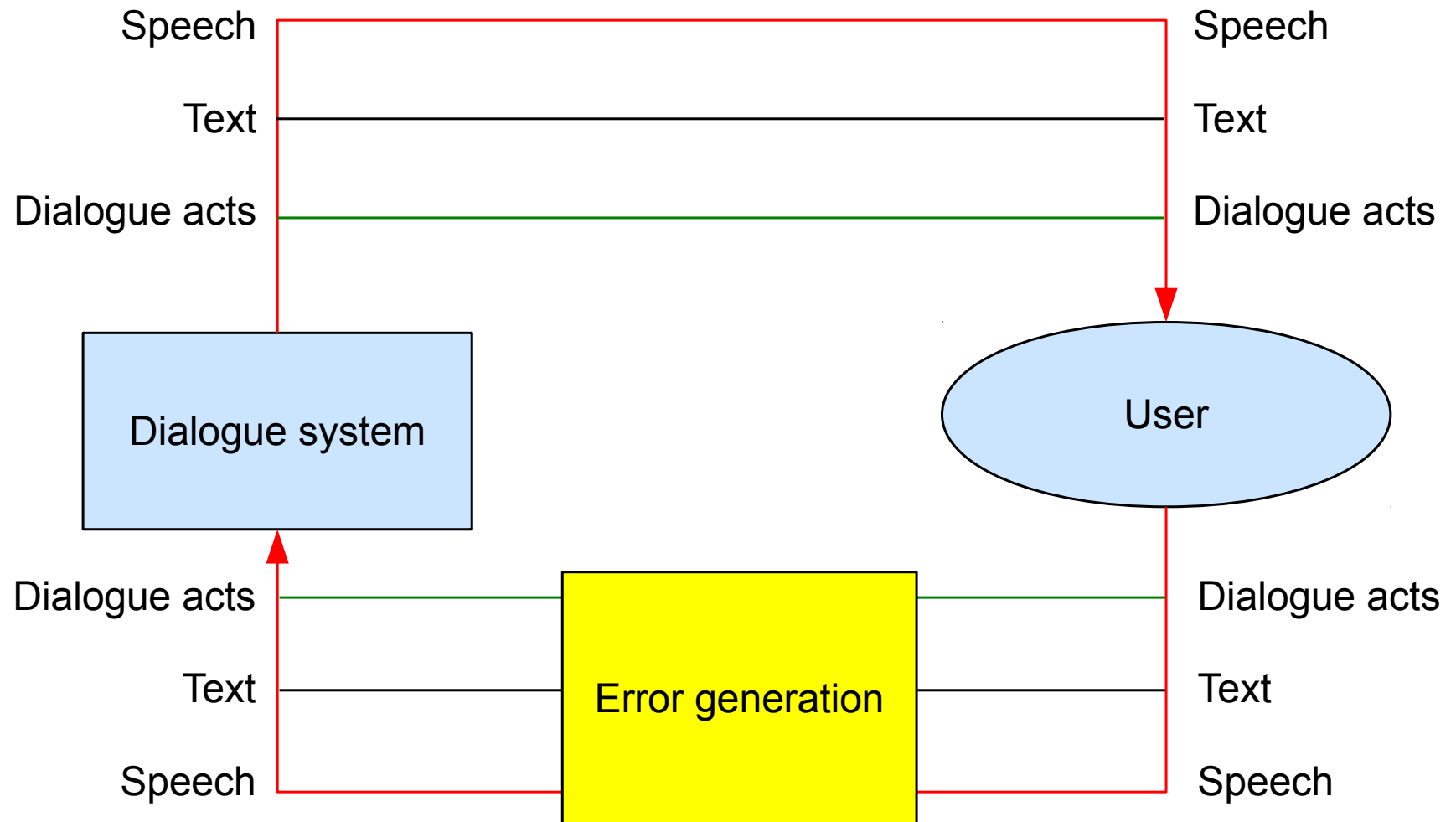


# Error simulation



# Types of user simulation

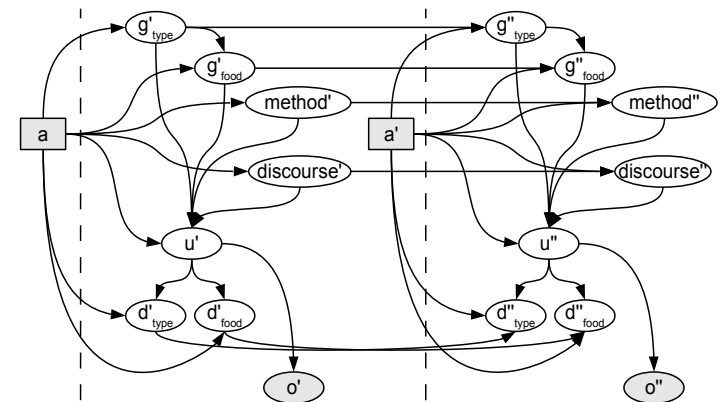
- Dialogue act level, Word level, **Speech level**



# Dialogue act level simulation

- Dialogue act level
  - S: request(food\_type)
  - U: inform(food\_type=Chinese)
- Many different implementations
  - bigram model for dialogue act types and random sampling slots
  - agenda based simulator
  - using user model to sample user dialogue acts

- Output dialogue act confusion model





# N-gram level user simulation

- Dialogue act level simulation
- Predict the next dialogue act given some context

$$P(d|\dots)$$

- Sample from a distribution
  - sampling insures variability in the output
- Typical context:
  - None – unigram models
  - Previous dialogue act – bigram models

# N-gram dialogue act user simulation

- In general,  $P(d|\dots)$  is too complex to model explicitly
- Model dialogue act type independently of the slots

$$P(dat|\dots) \approx P(dat_{t+1}|dat_t) \approx P(dat_{t+1})$$

request → inform

inform → inform

inform → deny

inform → confirm

....

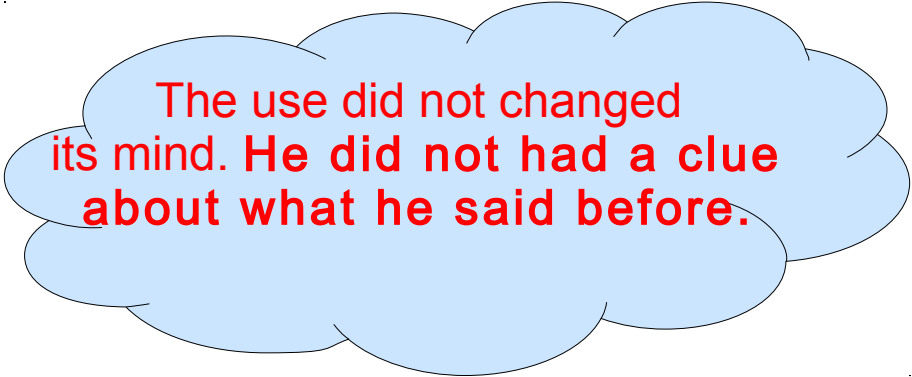
- Slots names depends on DAT

$$P(sn_{t+1}|dat_{t+1})$$

- or are drawn from a uniform distribution

# N-gram dialogue act user simulation

- Can be easily learned from the collected data
- These users are not rational due to small context considered
- Values are very often randomly selected
  - there is not enough data to estimate probabilities for this
- Example:
  - S: request(pricerange)
  - U: inform(pricerange=**cheap**)
  - S: request(pricerange)
  - U: inform(pricerange=**expensive**)



The user did not change its mind. He did not have a clue about what he said before.

# Agenda based simulations

- Inspired by agenda based approach for SDS:
  - X Wei and Al Rudnicky. 1999. An agenda-based dialog management architecture for spoken language systems. In Proc. of IEEE ASRU. Seattle,WA.
- In this case, the dialogue state is factored into
  - the goal – G
  - the agenda – A
- The goal ensures that the user behaves in a consistent, goal-directed manner.
- For example in tourist information domain, the goal can be further factored into
  - requirements - R
  - constraints - C

# Agenda based simulations

- The user agenda is a stack-like structure containing the pending user dialogue acts
- Note: the agenda is a special way of representing a dialogue policy and the context history at the same time
- At the start of the dialogue the agenda contains
  - Inform acts about all constraints
  - Request act for all requirements
  - Bye dialogue act

# Agenda update

- Based on the systems responses, update the agenda of the simulator, e.g. what to say next

$$P(a_{t+1} | a_t, g_t, u_{a:t})$$

- Although it can have probabilistic representation, it is mostly done deterministically, e.g.:
  - If system misunderstands then the user simulator puts on the agenda a correction
- Since agenda operates with a stack-like structure, the allowed operation are:
  - push and pop acts to and from agenda
  - the dialogue end when *bye()* act is popped of the stack

# Example

$C_0$	=	$\left[ \begin{array}{l} type = bar \\ drinks = beer \\ area = central \end{array} \right]$	Usr 2	No, beer please!
$R_0$	=	$\left[ \begin{array}{l} name = \\ addr = \\ phone = \end{array} \right]$	Sys 2	A bar serving beer, correct?
Sys 0		Hello, how may I help you?	$A_3$	= $\left[ \begin{array}{l} affirm() \\ inform(pricerange = cheap) \\ inform(area = central) \\ request(name) \\ request(addr) \\ request(phone) \\ bye() \end{array} \right]$
$A_1$	=	$\left[ \begin{array}{l} inform(type = bar) \\ inform(drinks = beer) \\ inform(area = central) \\ request(name) \\ request(addr) \\ request(phone) \\ bye() \end{array} \right]$	Usr 3	Yeah something cheap in the town centre.
Usr 1		I'm looking for a nice bar serving beer.	Sys 3	Murphy's on Main Square serves cheap beer.
Sys 1		Ok, a wine bar. What pricerange?	$A_4$	= $\left[ \begin{array}{l} request(phone) \\ bye() \end{array} \right]$
$A_2$	=	$\left[ \begin{array}{l} \underline{negate(drinks = beer)} \\ inform(pricerange = cheap) \\ inform(area = central) \\ request(name) \\ request(addr) \\ request(phone) \\ bye() \end{array} \right]$	Usr 4	Ok, and what's the phone number?
			Sys 4	The number is 796 69 94.
			$A_5$	= $\left[ \begin{array}{l} bye() \end{array} \right]$
			Usr 5	Thanks, goodbye!

# Agenda update

- Agenda update tends to be complex
  - due to limited number of ways of updating the agenda
  - **!!! the update is limited to make its use easier !!!**
- Agenda update is unnecessary factored into too small and specific operations
- Some trivial things are just too complex to be achieved with push and pop operations



# ISU based approach

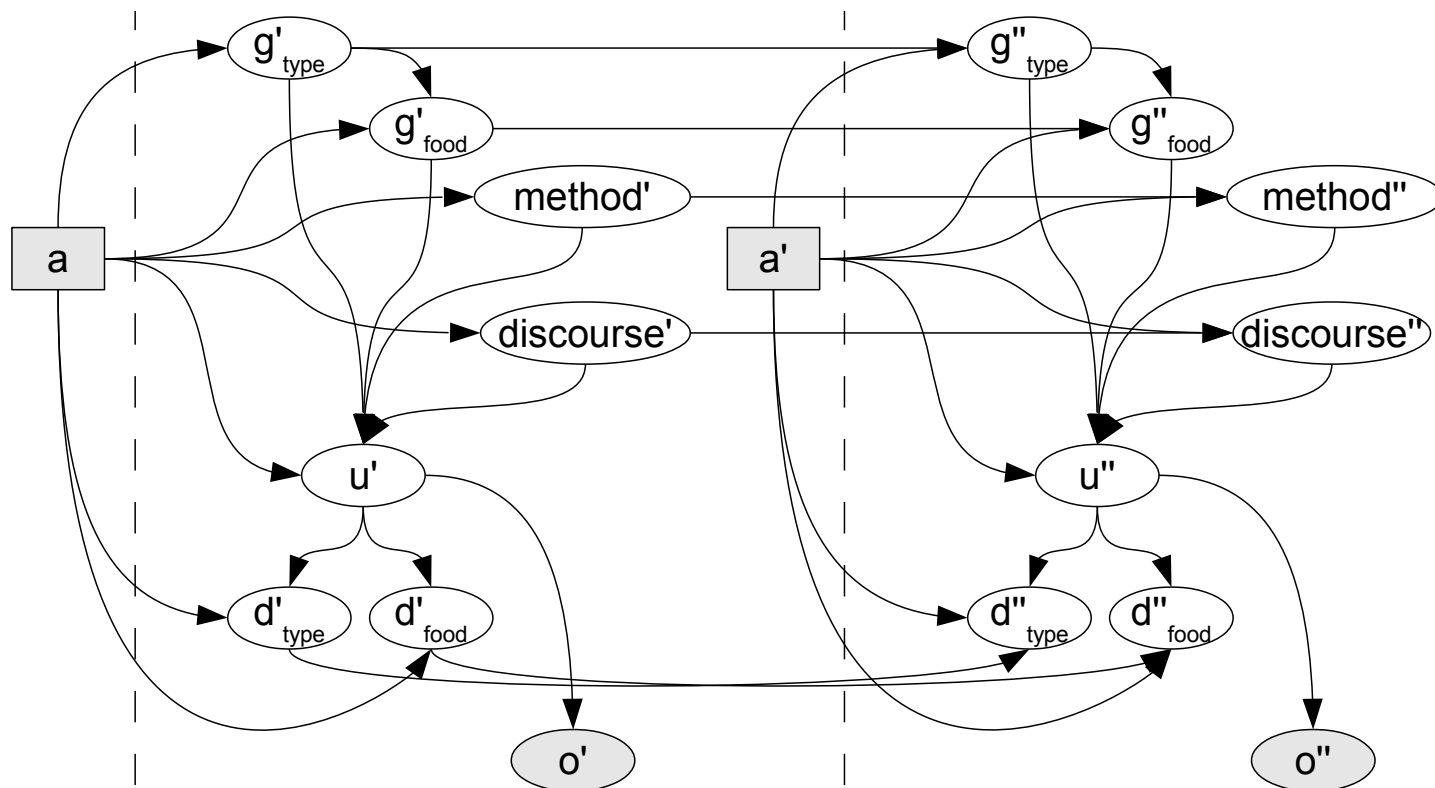
- ISU based approach would simply build a dialogue system behaving like a user:
- Such system:
  - would track what the dialogue systems said
  - compared that information with the goals of the user
  - then made it decisions using either:
    - a handcrafted policy
    - a stochastic policy
- To limit complexity of the dialogue policy
  - summary action can be used
  - they can be tuned to needs of a user simulator

# ISU based approach

- It is like building a dialogue system without reinforcement learning
  - we can learn stochastic policy using maximum likelihood
- We do not want to be better than a user
- We want to do exactly what a user is doing
- Though, we could use reward function to constrain our model to, for example, cooperative policies
  - e.g.: a user wants to succeed in a dialogue

# Bayesian approach

- Use model similar to the dialogue model in a dialogue manager



# Sample from the Bayesian model

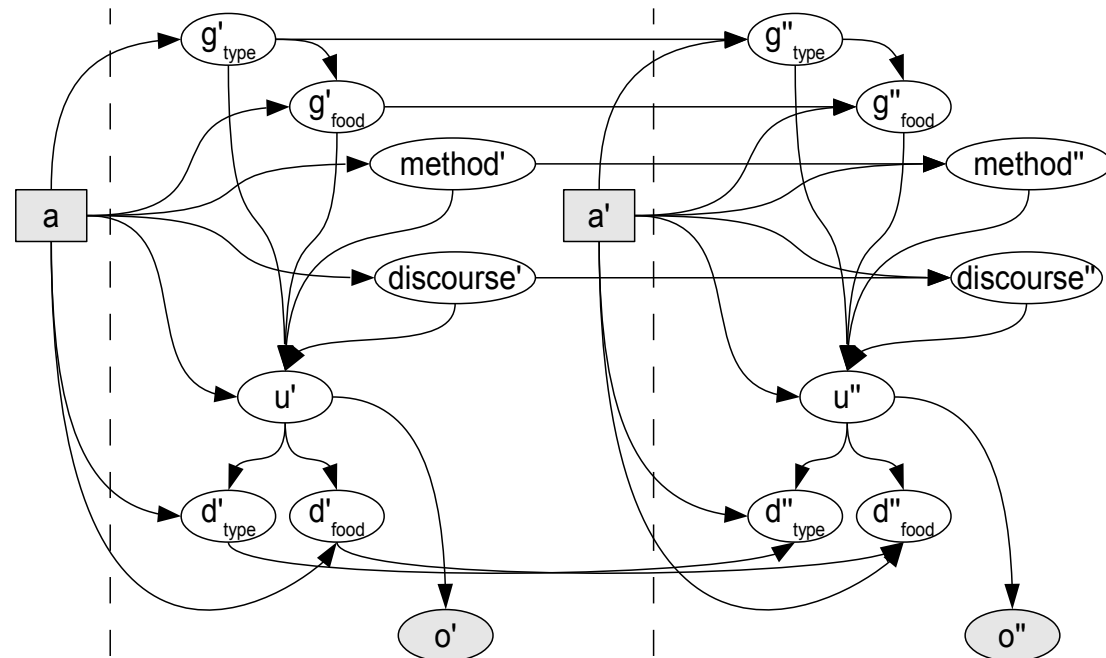
- Sample next dialogue state

$$s \sim P(s_{t+1} | s_t, a_t)$$

- Sample next user action

$$o_{t+1} \sim P(s_{t+1})$$

- Problem
  - it is not always rational
  - this is similar to the N-gram dialogue act level simulations



# Bayesian approach

- Ideally the model for simulations and the dialogue system would be the same
- However, as typical for generative models:
  - generative models can be useful for classification
  - without being good for generation!
- The model for user simulations must be:
  - more constrained → more (probabilistic) handcrafting
  - it needs longer context
  - it turns out to be complex and hard to maintain
- However, it appears to be the principled way to do

# Error generation

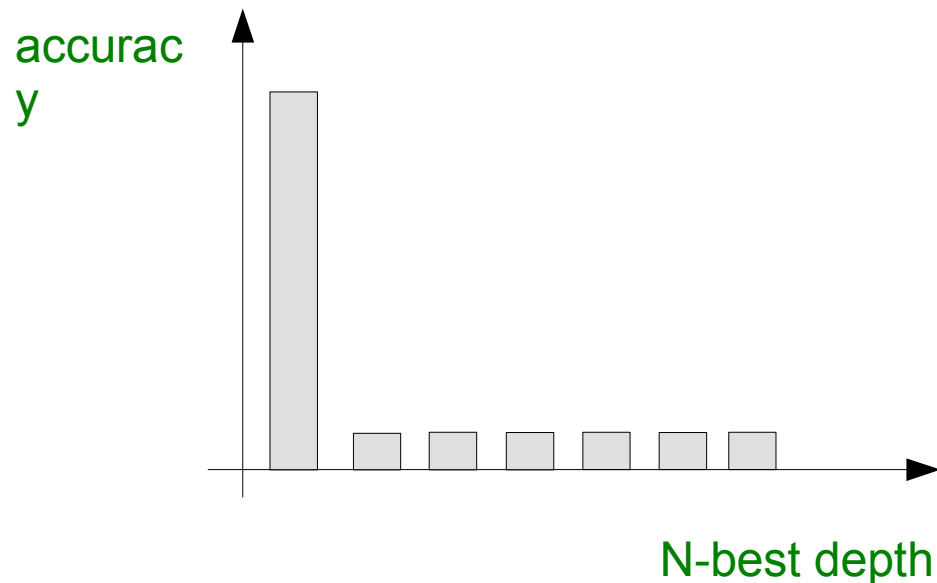
- So far, we wanted from the user simulator reasonable behaviour
- However, for training a robust dialogue system we need a model of errors: ASR + SLU

inform(food=Italian)&request(price): 1.0	→	inform(food=Italian)	: 0.50
	→	inform(food=Italian)&request(price)	: 0.20
	→	request(price)	: 0.20
	→	confirm(food=Italian)	: 0.05
	→	null()	: 0.05

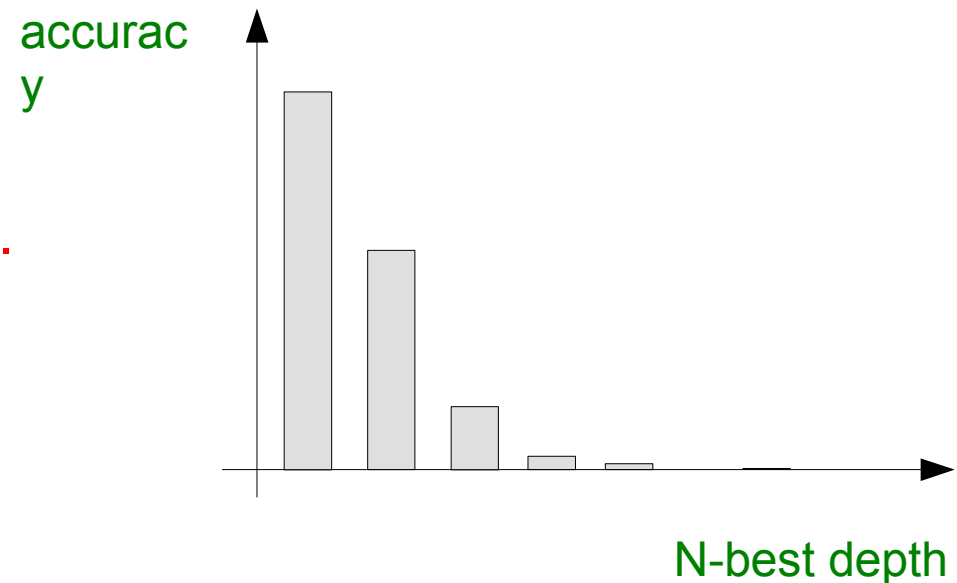
- Ideally, we can control this error model using a simple variable: error rate
  - measure of the number of errors in the N-best lists / top hypothesis

# Error generation

- We want to control both the accuracy of the top hypothesis and the distribution of the accuracy at different positions in the N-best list

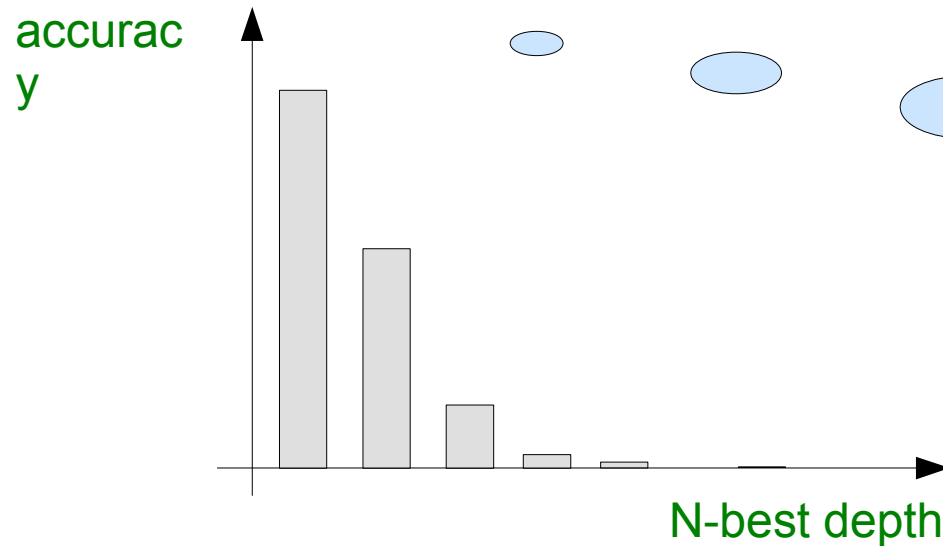


VS.

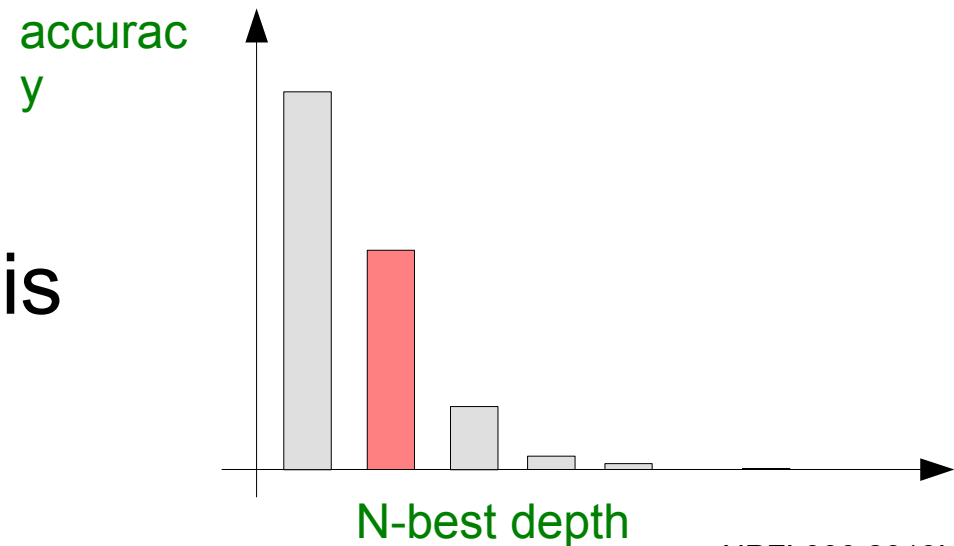


# Model for error generation

- Error model : accuracy =  $P(\text{position in the N-best list})$



- Sample position of the correct hypothesis

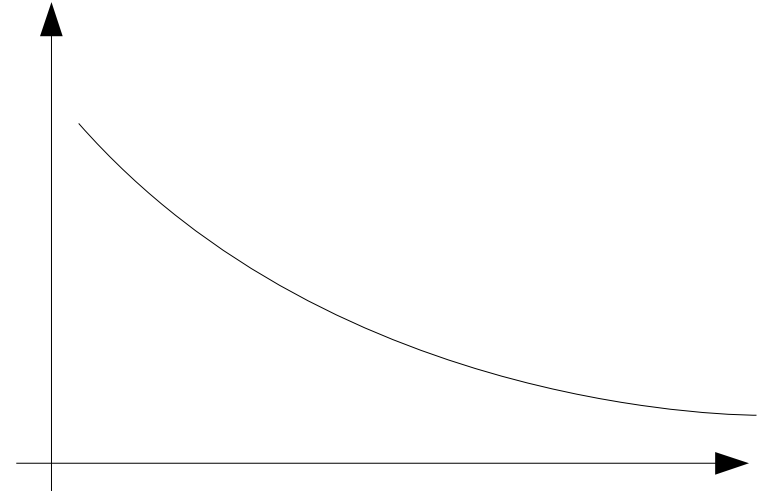




# Model for error generation

- Example implementation

$$P(n|err; \alpha) \approx e^{-\alpha \frac{n}{err}}$$



- Having a position of the correct hypothesis

inform(food=Italian)&request(price): 1.0	→	...	: 0.50
	→	inform(food=Italian)&request(price)	: 0.20
	→	...	: 0.20
	→	...	: 0.05
	→	...	: 0.05

- We have to fill in the missing hypotheses

# Model for error generation

- Sample from a dialogue act confusion model

$$d \sim P(d_{err} | d, err, n)$$

- This is usually factored into

$$d \sim \prod_i P(dact_{err}^i | dact^i, err, n) P(sn_{err}^i | dact_{err}^i, sn^i, err, n) \\ P(sv_{err}^i | sn_{err}^i, err, n)$$

- And then the individual probabilities are usually handcrafted:
  - typically ignores **err** and **n**

# Trained model for error generation

- Use logistic (multiclass) regression to train the model from errors generated by ASR, SLU

$$P(dact_{err} | dact, err, n) \approx e^{\theta^T \Phi(dact_{err}, dact, err, n)}$$

$$P(sn_{err}^i | act_{err}^i, sn^i, err, n) \approx e^{\theta^T \Phi(sn_{err}^i, act_{err}^i, sn^i, err, n)}$$

$$P(sv_{err}^i | sn_{err}^i, err, n) \approx e^{\theta^T \Phi(sv_{err}^i, sn_{err}^i, err, n)}$$

- Every time you change ASR (AM, LM) or SLU, you have to re-train these models

# Word level user simulation

- In addition to dialogue act level simulation
  - implements reasoning and goal oriented behaviour
- It include the generation of sentences
  - e.g. using a corpus of dialogue act annotated sentences
- The confusion model then operate on words
  - it learns mapping between correct word strings and confused words
  - it can use again pre-recorded and consequently recognised corpus
  - however, some generalisation is needed to obtain larger variability in the output

# Word level user simulation

- Tries to solve the problem of similarity of
  - `inform(type=bar)` - “A bar please!”
  - `inform(drinks=beer)` - “Uh beer please!”
- and dissimilarity of
  - `inform(type=restaurant)` - “A restaurant please!”
- In general, semantically similar items does not have to be similar on the word/phone level

# Word level simulation

- Phone level simulation would be ideal
- However, this is far too complex
- It is more efficient to learn how to confuse individual words or word sequences
- Note: this is still too computationally expensive anyway

# Word level confusion model

$$\begin{aligned} P(\tilde{a}_u | a_u) &= \sum_{\tilde{w}_u} \sum_{w_u} P(\tilde{a}_u, \tilde{w}_u, w_u | a_u) \\ &= \sum_{\tilde{w}_u} \underbrace{P(\tilde{a}_u | \tilde{w}_u)}_{\text{semantic decoder}} \sum_{w_u} \underbrace{P(\tilde{w}_u | w_u)}_{\text{confusion model}} \underbrace{P(w_u | a_u)}_{\text{utterance generation}} \end{aligned}$$

$$P(\tilde{w}_u | w_u) = \sum_{\lambda} P(\tilde{w}_u, \lambda | w_u)$$

# Utterance generation model

- For error generation, a simple template generation model can be considered

Source:           inform(food=Chinese, pricerange=cheap)  
CHINESE FOOD IN THE CHEAP PRICERANGE

Template:        inform(food=\$X, pricerange=\$Y)  
\$X FOOD IN THE \$Y PRICERANGE

Unseen:          inform(food=French, pricerange=expensive)  
FRENCH FOOD IN THE EXPENSIVE PRICERANGE

- Templates can be easily derived from the corpus using a database of slot values
- $P(\text{template}|\text{dialogue act})$  can be also collected



# Word confusion model

- At the word-level, ASR confusions can be viewed as translations of a source utterance  $w_s$  to a confused target utterance  $w_t$

$$P(\tilde{w}_u | w_u) = \sum_{\lambda} P(\tilde{w}_u, \lambda | w_u)$$

I	WANT	AN	EXPENSIVE	HOTEL	PLEASE
1	2	3	3	4	5
			/		
1	2	3		4	5
	ONE	INEXPENSIVE		HOTEL	PLEASE

# Word confusion model

- There are many ways how to define this confusion model
- Nevertheless, all depend on some alignment of between  $w_s$  and  $w_t$
- The easiest is to use a Levenshtein distance
  - automatically generated alignment
  - word level with costs defined on the letter/phone level
- This is a bit different from machine translation since we work with the same language on both sides

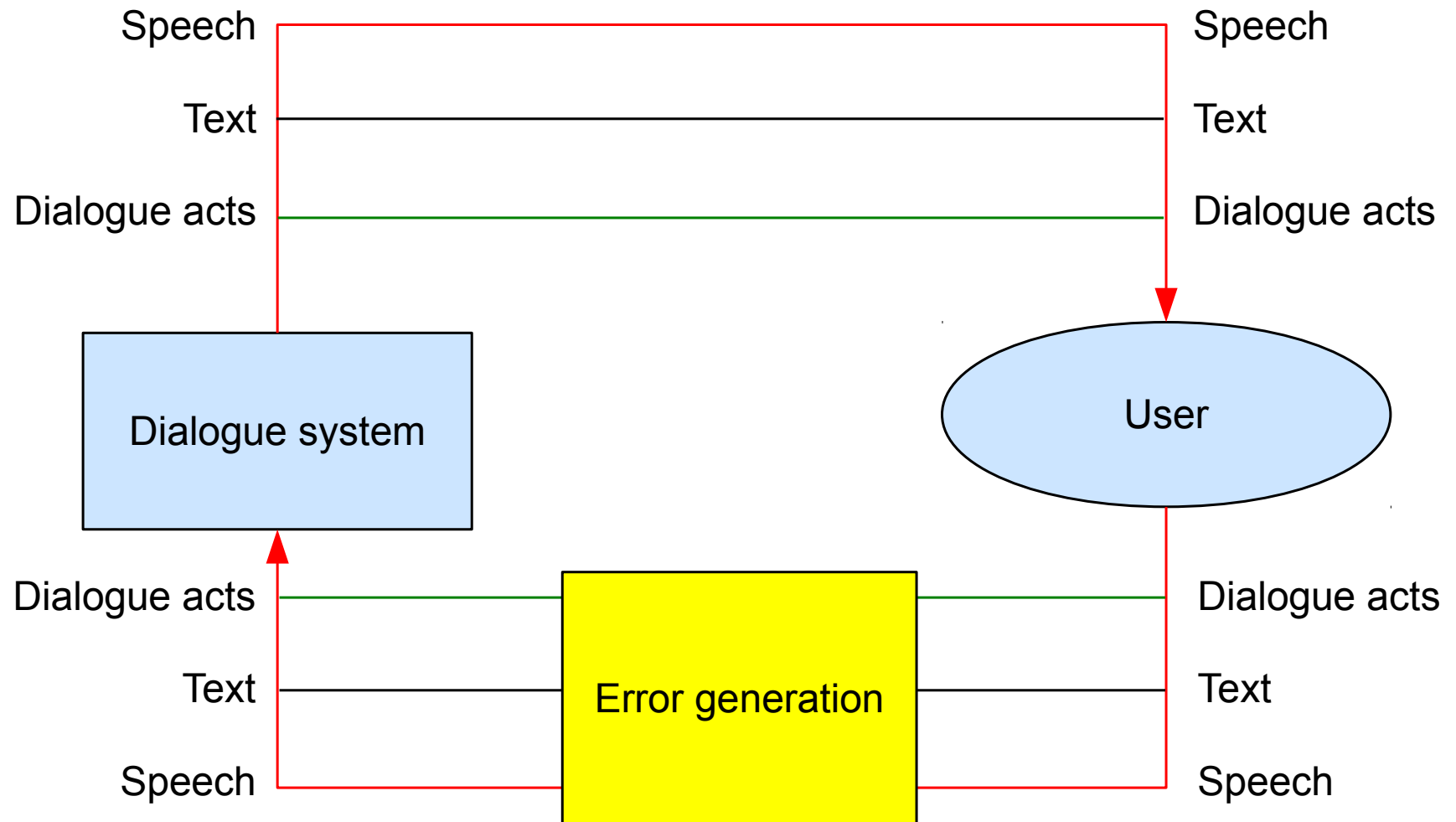
# Word confusion model

- Once you get the alignments,
  - you can compute word / word sequence confusion model
  - e.g. "A BAR" maps to
    - "ALL", "ART", "A BAR", "A BAR", "A CAR", "BAR", "BAR", "BAR", "BAR", "CAR"
  - and corresponding probabilities
- Then the following can be used compute the probability of the alignment

$$\begin{aligned} P(\tilde{w}_u, \lambda | w_u) &\approx \prod_i P(\tilde{w}_{\lambda_i} | w_{\lambda_i}, w_{\lambda_{i-1}}, w_{\tilde{\lambda}_{i-1}}) \\ &\approx \prod_i P(\tilde{w}_{\lambda_i} | w_{\lambda_i}) \end{aligned}$$

# Types of user simulation

- Dialogue act level, Word level, **Speech level**



# Speech level user simulation

- In addition to word level simulation, it generate speech
- The generation of speech is done using TTS on correctly generated text
- The noise is added by simply mixing generated speech and some artificial noise
- The noise is controlled by its volume and type
- Problems:
  - The generated text has low variability and can be unnatural in some cases
  - The synthesised speech has low variability
  - Obviously, the acoustic noise is not responsible for all uncertainty in the SDS input

# Thank you!

Filip Jurčiček

Institute of Formal and Applied Linguistics  
Charles University in Prague  
Czech Republic

Home page: <http://ufal.mff.cuni.cz/~jurcicek>

