# NPFL099 - Statistical dialogue systems

# Dialogue system evaluation

Filip Jurčíček

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

Home page: http://ufal.mff.cuni.cz/~jurcicek
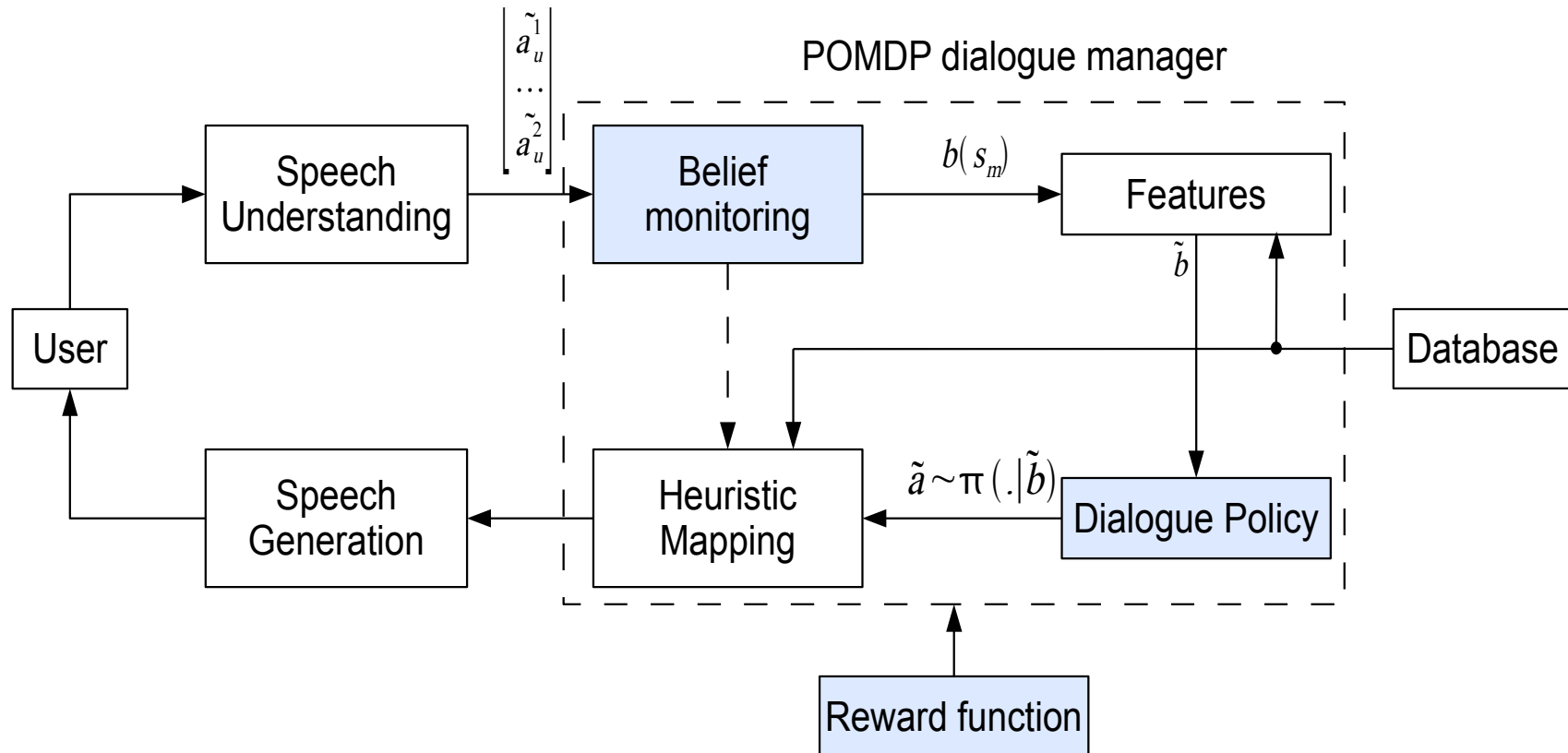
Version: 14/05/2013

# Outline

- Evaluation of dialogue systems

- Laboratory evaluation

- Crowdsouricing

- Real user evaluation

# Evaluating of SDS

- evaluate each component separately

- or in the context of all others

# Metrics - subjective

- Typically a feedback form

  - Did you find all the information you were looking for?
    - to evaluate the dialogue manager

  - The system understood me well
    - to evaluate the spoken language understanding component

  - The phrasing of the system's responses was good
    - to evaluate the language generation component

  - The system's voice was of good quality
    - to evaluate the speech synthesizer

# Metrics - subjective

- You do not want to ask many questions

  - It is boring

  - The answers typically correlate

- Ideally yes / no questions

- or selection from N options
  - Likert scale

  - ``strongly disagree'', ``disagree'', ``lightly disagree'', ``slightly agree'', ``agree'', and ``strongly agree''.

# Metrics - objective

- So far, we talked about subjective metrics

- Some times
  - subjective metrics are not available

  - objective (automatic) metrics can be better ???

- Objective metrics
  - PERSEVAL – Walker et. al
    - trainable model from a corpus of human ratings
    - most explanatory features are accuracy of ASR, the length of the dialogue

  - BETTER: was the call routed to a human operator?
    - not always applicable

# Evaluation in a laboratory

- In controlled environment
  - noiseless or with generated back ground noise

- Typically
  - each user gets training
  - is supervised by an assistant
  - users rating is controlled

- When interacting with the SDS
  - user is given a goal
  - the assistant could point out
    - errors in rating
    - missing constraints in the goal

# Evaluation in a laboratory

- It is time consuming

- Search for subjects among colleagues or students

- You have to make appointments
  - some people does not show up

- Expensive – in CAM, we paid £15 for an hour

- Still, we could not get enough subjects

# Evaluation using crowdsourcing

- Similar setup but hiring users differently

- Amazon Mechanical Turk users
  - toll-free phone number in USA

  - mostly native English speakers from USA

  - some Canadians

  - many Indians

  - some non native speakers of English from USA

# Evaluation using crowdsourcing

- Instead of coming to a lab, subjects were presented with a web page

- Web interface
  - To instruct users
  - To give tasks
  - To collect feedback

- Phones used to deliver voice
  - Calls routed using SIP to Cambridge, UK

# Web interface

# Evaluation using crowdsourcing

- Relatively easy to get users
  - between 100 – 200 calls a day

  - better to ask TURKs to test a system than to ask colleagues ;-)

- Cheap – minimum wage
  - we paid  $6 for an hour

- Toll-free phone number cheap
  - $0.02 per minute

# Evaluation with real user

- TURKS are not real user
  - they are still paid
  - their rating is in some extent random
    - though this is true for all humans
    - unless you go to the recommended restaurant, it is hard to rate usefulness of the SDS recommendation

- Would be better to have real users interested in using the SDS
  - only some have such applications
  - e.g. Speech Cycle, Nuance, France Telecom
    - have tens of thousand calls a month

# Evaluation with real user

- Still, the rating does not have to be consistent

- The reward can be delayed

- I will know that the appointment booking was successful only when the technician comes on the date I wanted

  - FT: appointment booking application

- You do not want ask all users
  - therefore automatic metrics are preferred

# Metrics – how many user do we need?

- Many!!!

- Imagine testing a system

  - the success rate is about 50 – 60 percent
  - when you collect 500 dialogues then the 95% confidence interval +- 5%

- Using parametric tests, a difference of less than 5% is not statistically significant

# Example: MTURK trial

- Amazon Mechanical Turk users

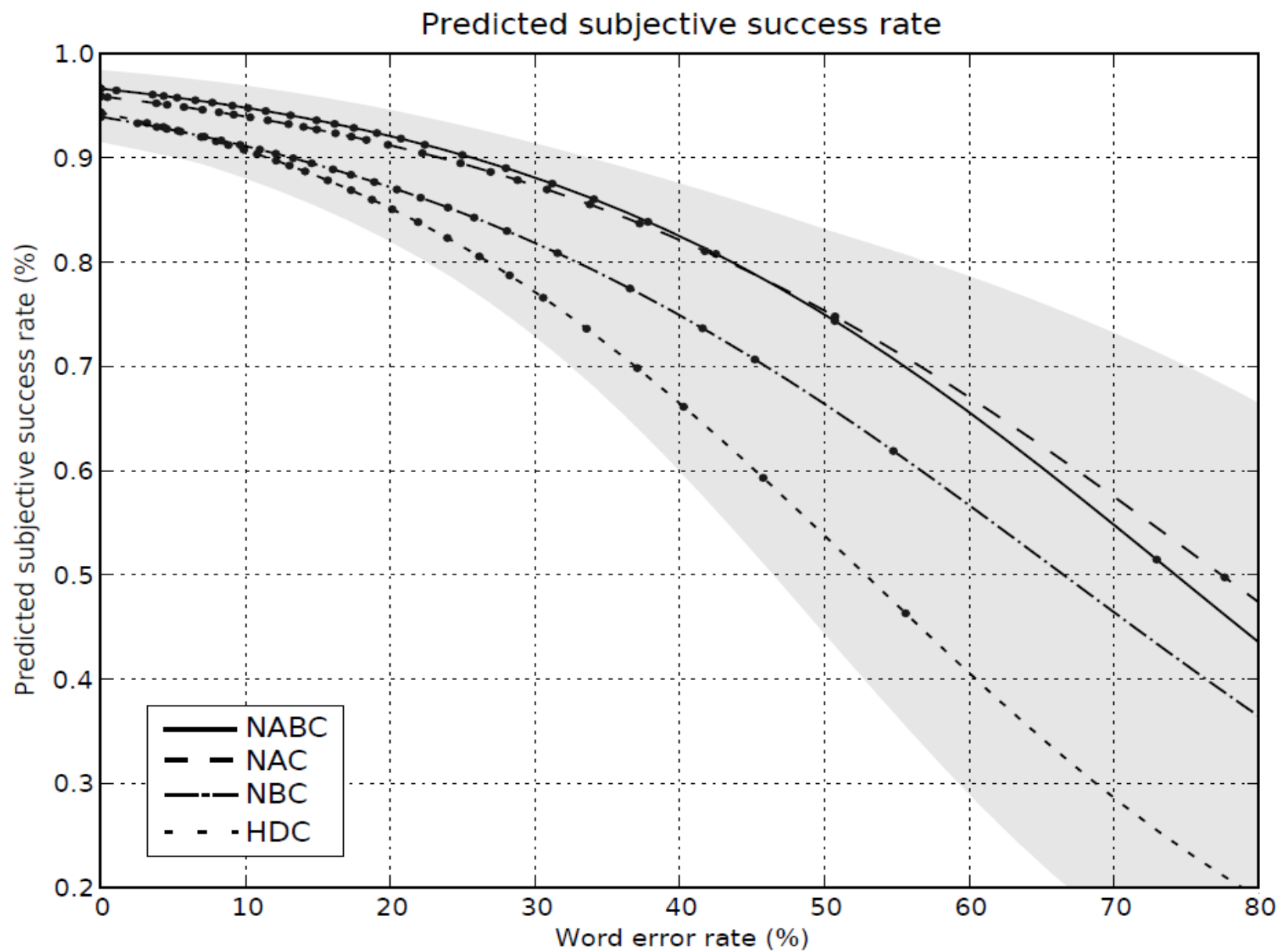| | |
|---|---|
| The number of calls | 2354 |
| The number of turns | 25289 |
| The number of users | 164 |
| ASR Word error rate | 20.1% |
| Length of the audio | 70 hours |
| Average length of a call | 1:47 min |

# Results: MTURK trial

- ## Metrics
  - Subjective success rate – user ratings
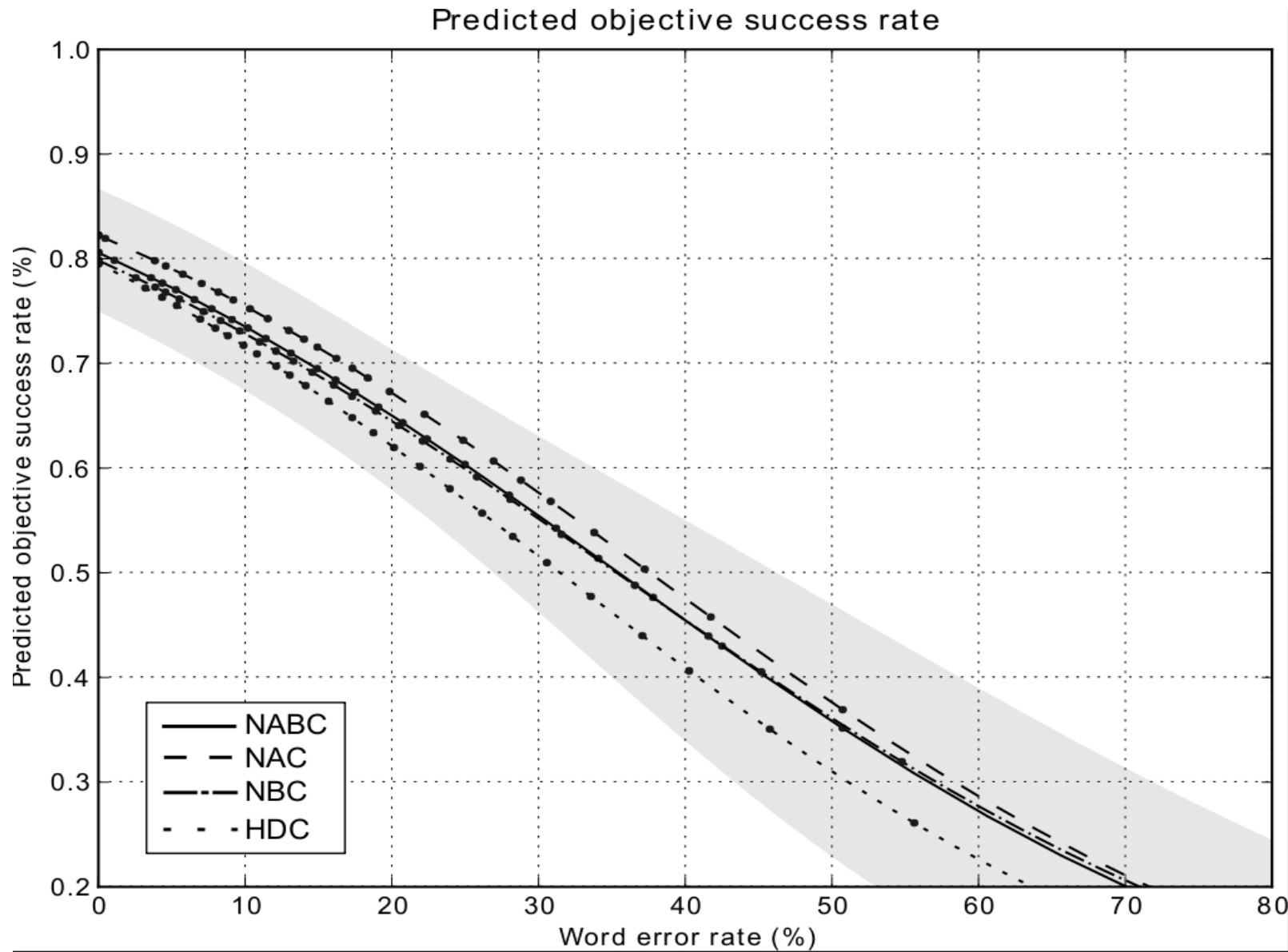  - Objective success rate – automatically derived

| System | # calls | Subjective Success Rate | Objective Success Rate |
|---|---|---|---|
| HDC | 627 | 82.30% ($\pm$2.99) | 62.36% ($\pm$3.81) |
| NBC | 573 | 84.47% ($\pm$2.97) | 63.53% ($\pm$3.95) |
| NAC | 588 | 89.63% ($\pm$2.46) | 66.84% ($\pm$3.79) |
| NABC | 566 | 90.28% ($\pm$2.44) | 65.55% ($\pm$3.91) |

- ## This does not say much about the performance at different error rates

# Results: subjective scores



Predicted subjective success rate

# Results: objective scores

# Thank you!

Filip Jurčíček

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

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