

End-to-end Neural Dialogue Systems

Ondřej Dušek

VOCALLS AI Afternoon

19.10.2022



Charles University
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics

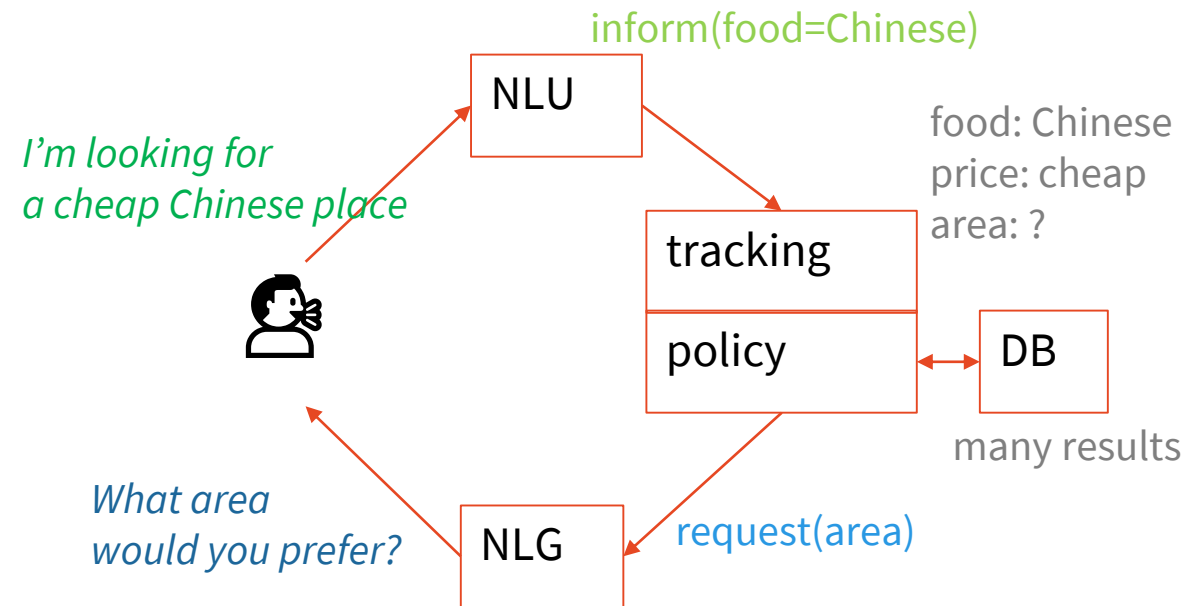


unless otherwise stated

Dialogue systems

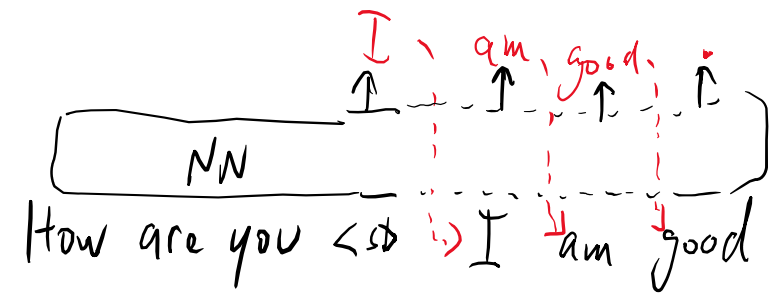
Standard architecture (text-based):

- **Natural language understanding**
 - user utterance → user dialogue act (DA)
- **Dialogue state tracking**
 - update user preferences based on DA
- **Dialogue policy**
 - choose next action based on state
 - system DA
 - consult DB if needed
- **Natural language generation**
 - express system DA as a sentence



End-to-end neural systems

- Standard approach: **Pipeline**
 - each module built separately
 - NLU typically trained (or combined with rules)
 - Tracking could be either
 - Policy, DB, NLG: rules, templates
 - problem: error accumulation, costly to implement
- **End-to-end models:** remove the pipeline
 - NLU/Tracking/Policy/NLG is one neural network
 - generate response word-by-word
 - DB must be external → 2-step operation
 - trained on (many) example dialogues



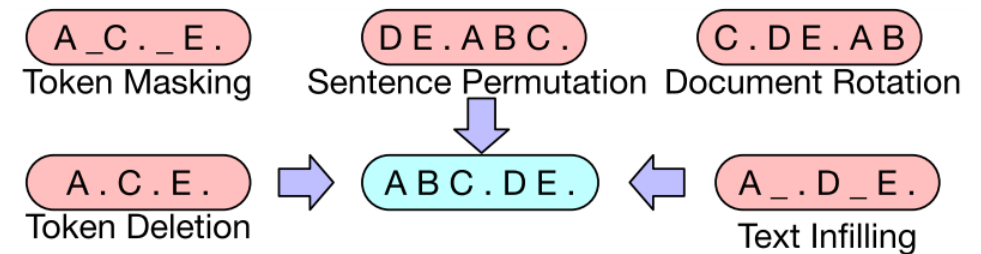
Neural language models

- **Transformer** neural architecture

- (sub)word representation: **embedding** = vector of numbers
- **blocks: attention** (combining context) + **fully-connected** (abstracting)
- predicting next (sub)word = classification: choosing 1 out of ca. 50k (probabilistic)
- trained from data: initialize randomly & iteratively improve

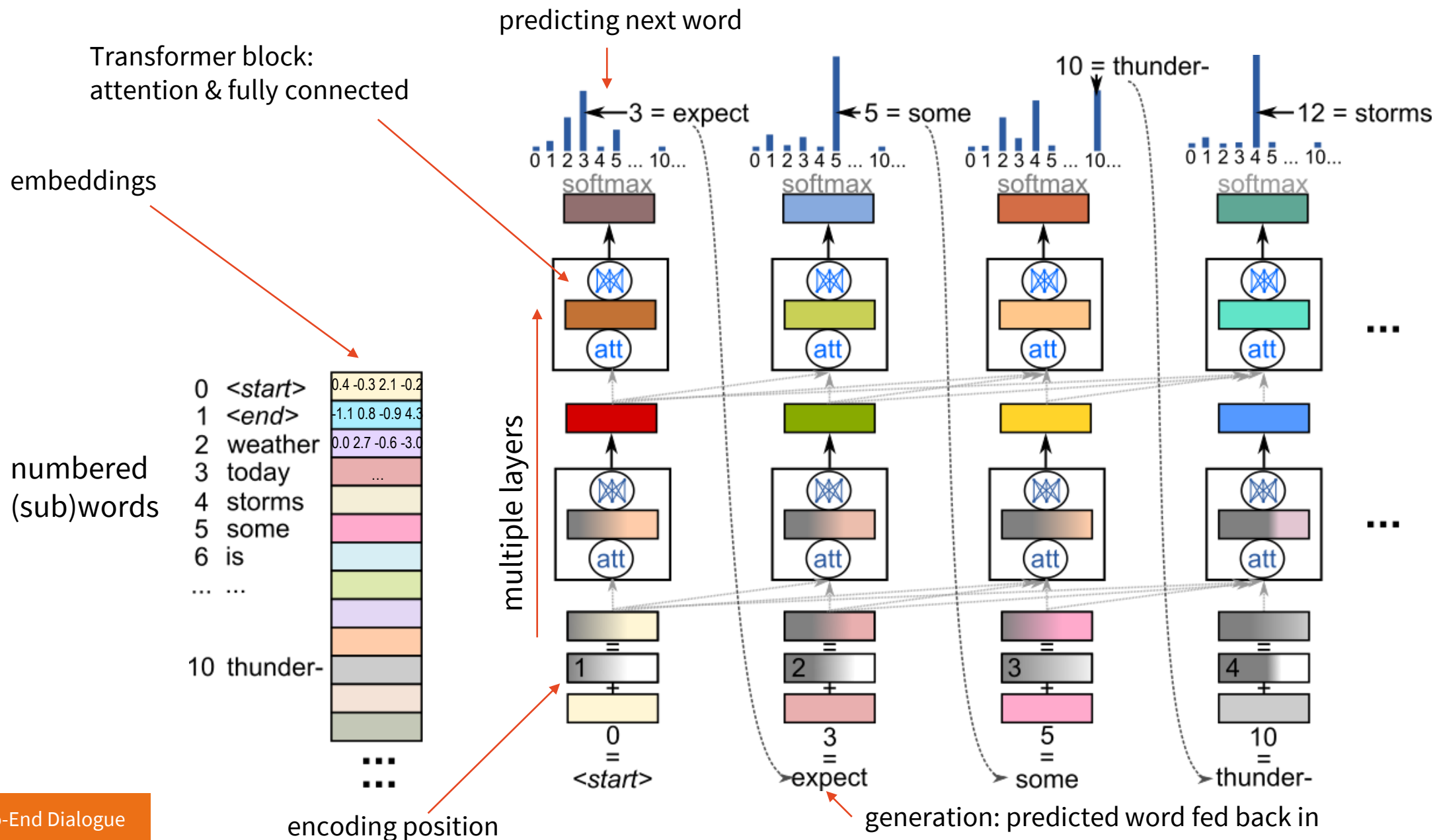
- **Pretrained models**

- Transformers trained on vast amounts of data
- Self-supervised training: just naturally occurring text & simple tasks
 - predicting next word
 - predicting masked word
 - fixing corrupt sentences
 - ...
- Lot of them released online, plug-and-play



<https://github.com/huggingface/transformers>

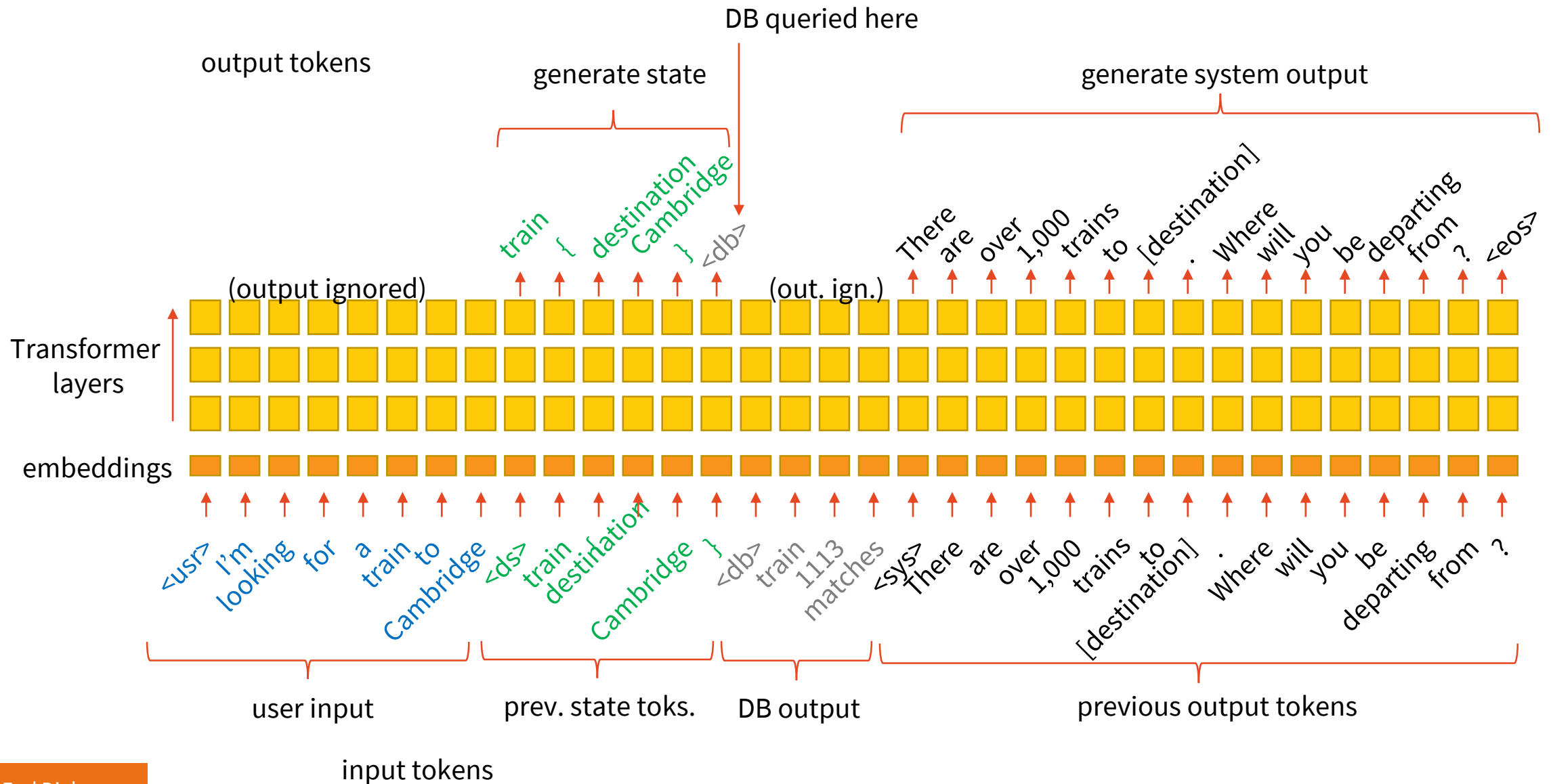
Transformer neural language model



End-to-end Neural Dialogue with GPT-2

- **GPT-2**: one of most popular pretrained language models
 - Transformer <https://huggingface.co/gpt2>
 - pretrained on next-word prediction
 - 8M docs, 40GB data from the web
- **Dialogue**: GPT-2 finetuned (=further trained) on dialogue data (MultiWOZ)
- Single neural network, fits Transformer operation <https://github.com/ufal/augpt>
- Multi-step, all word-by-word:
 1. feed in dialogue context (← ignore generation outputs for this bit)
 2. generate dialogue state (as text)
 3. query DB
 4. feed in DB results as text (← ignore outputs)
 5. generate response

End-to-end Neural Dialogue with GPT-2



Problems & solutions

- **Needs a lot of data** & annotation (1000s of dialogues)
 - costly, may be noisy
 - → transfer learning, data augmentation
- **Hallucinates** sometimes
 - may generate factually incorrect outputs, hard to control
 - → data cleaning, consistency training (corrupt data & train to detect this)
- **Repetitive/dull** outputs
 - settles for the most frequent output
 - → sampling (randomness)
- Still a long way to go
 - ~70% correct/successful dialogues
 - still needs a lot of data

Thanks

Contact me:

odusek@ufal.mff.cuni.cz
<http://ufal.cz/ondrej-dusek>

Credits:



Vojtěch Hudeček



Jonáš Kulhánek



Tomáš Nekvinda

Our work:

<https://github.com/ufal/augpt>

These slides:

<https://bit.ly/e2e-ds-vocalls>