NPFL099 Statistical Dialogue Systems **11. Multimodal Systems** (+some notes on domain adaptation)

http://ufal.cz/npfl099

Ondřej Dušek, Vojtěch Hudeček & Zdeněk Kasner 12. 12. 2022



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



Multimodal Dialogue Systems

- adding more modalities to voice/text
 - input:
 - touch
 - drawing
 - gaze, gestures, facial expressions
 - voice pitch/tone
 - image
 - output:
 - graphics
 - gaze, gestures, facial expressions, body movement
- either traditional/modular and mostly rule-based systems, or very experimental (not much use in practice)

Standard Multimodal DS Schema

- basically the same as voice/text DSs
- adding multiple input modules
 - for multiple modalities
 - each with its own NLU-like interpretation
 - interpretations are merged
- multiple output modules
 - each with its own generation
 - dialogue manager output is split
- typically ready-made off-the-shelf modules
 - it's too complex/costly to build these custom



Smart Devices

- Phones, wearables, smart speakers with a display
 - incl. Google Assistant, Alexa & Siri
 - admittedly not so much dialogue, more of commands
 - cloud-based operation for most
- Input
 - touch: active & passive gestures (touch/accelerometer)
 - "raise to speak"
 - rarely visually sensing gestures
 - doesn't support gaze
- Output
 - graphics: card interface
 - generation functions rule-based/low-level

https://www.wareable.com/android-wear/how-to-use-voice-commands-on-android-wear https://www.cnet.com/reviews/amazon-echo-spot-review/ https://www.macrumors.com/how-to/use-siri-raise-to-speak-watchos-5/



"Classical" Multimodal Systems

- closed-domain task-oriented dialogue systems
- map-based: town information with map input & output
 - touch / pen drawing, map display
 - reacting to zooming, area selection
 - handwriting recognition (as alternative input)
 - similar to Google Assistant, but more interactive
- in-car: voice & button control
- custom architectures
 - off-the-shelf modules
 - rule-based touch input processing



- S: I found 3 albums by The Beatles in your collection <shows listing on screen>
- U: Play the third one.
- U: Which songs are on this one?
 - <selects an album from listing on screen>

Virtual Agents

https://youtu.be/ejczMs6b1Q4 https://vhtoolkit.ict.usc.edu/

- character face/full body
 - on screen or 3D projected (FurHat)
- a lot more outputs
 - full motion video facial expressions, gaze, gestures, body movement
 - a lot of it "automatic", designed to look natural/match what's said
- additional inputs gaze & facial expression
 - checking user engagement/sentiment
- dialogue management mostly rule-based
 - retrieval with non-linguistic inputs (Virtual Humans/SimSensei)
 - limited-domain custom rules (FurHat)
- tutoring/training, healthcare

FurHat





Virtual Humans SimSensei

(Al Moubayed et al., 2012) (Rushforth et al., 2009) (DeVault et al., 2014) https://doi.org/10.1007/978-3-642-34584-5_9 https://doi.org/10.1007/978-3-642-04380-2_82 https://dl.acm.org/doi/10.5555/2615731.2617415

Robots

- similar to virtual agents, but with actual hardware
 - different user's perception
 - body gestures more prominent
 - touching the robot is possible
 - situated deployment need to track user engagement
 - is the user still talking to the robot?
 - hardware limitations
 - mostly no facial expr./gaze output, some sensors missing etc.
- off-the-shelf robots (Nao, Pepper)
 - built-in & additional sensors (e.g. Kinect)
 - custom rule-based gesture generation
 - controlled via a computer (not autonomous)
- "receptionist" directions, information





Multi-party Dialogue

- Relevant for both virtual agents & robots
 - supported by most previously mentioned projects
- How to handle multiple counterparts?
 - users or other robots/virtual agents
- gaze/engagement/speech detection
 - who's speaking/looking etc.
- rules for multiple counterparts
 - switching gaze to address them
 - here, 3D is better than 2D (otherwise gaze ambiguous)
 - telling one to wait for another
- customer service, information

https://youtu.be/oOp4XP ziMw http://www.danbohus.com/

(Foster et al., 2012) (Bohus et al., 2014)

http://dl.acm.org/citation.cfm?doid=2388676.2388680 https://dl.acm.org/doi/10.5555/2615731.2615835 (Skantze & Al Moubayed, 2012) https://doi.org/10.1145/2388676.2388698





S: Are you here looking for Zack?

Interaction 1 (Socially inappropriate)

Interaction 2 (Socially appropriate)

One person, A, approaches the bar a	and turns towards the bartender
Robot (to A): How can I help you?	Robot (to A): How can I help you?
A: A pint of cider, please.	A: A pint of cider, please.
A second person, B, approaches the	bar and turns towards the bartender
Robot (to B): How can I help you?	Robot (to B): One moment, please.
B: I'd like a pint of beer.	Robot: (Serves A)
Robot: (Serves B)	Robot (to B): Thanks for waiting.
Robot: (Serves A)	How can I help you?
	B: I'd like a pint of beer.
	Robot: (Serves B)

Specific uses

- Air traffic controller training radar as a modality
 - multiple agents/systems representing pilots
 - radar charting each agent's behavior
 - single ASR, many TTSs
 - varied accents
 - all rule-based
 - very limited domain
 - bearings, flight levels



End-to-end Multimodal

- recent, experimental
- enhancing end-to-end DS architectures with image input
 - no video input
 - no avatars, facial expressions, gestures etc.
 - not much graphics output either
- also using off-the-shelf components
 - especially for image recognition ready-made convolutional architectures
 - textual parts based on known architectures (HRED, MemNN etc.)
- mostly just end-to-end prediction
 - pretrained image recognition parts are kept fixed, no end-to-end training

Pretrained convolutional nets

- Data: ImageNet Challenge
 - >1M images, 1000 classes
 - just classify the object in the image
 - CNNs are way better than anything that came before them
- AlexNet 1st deep CNN
 - 5 conv layers, ReLU activations, max pooling & 3 dense layers
- VGGNet improvement
 - more layers, smaller CNN kernels (3x3, 2x2 pooling with stride 2)
 - reduces # of parameters, same function



https://towardsdatascience.com/the-w3h-of-alexnet-vggnet-resnet-and-inception-7baaaecccc96

(Krizhevsky et al., 2012)https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks(Simonyan & Zisserman, 2015)http://arxiv.org/abs/1409.1556

Pretrained CNNs

- ResNet residual networks
 - trying to simplify the mappings found by CNNs
 - with regular CNNs, deeper might not be better (vanishing gradient problem)
 - "shortcuts": adding identity / linear projection to convolutions
 - learning a residual CNN mapping ("what projection can't handle")
 - allows much deeper networks alleviates vanishing gradients
- Inception more CNN kernels in parallel
 - for detecting different-sized object features
 - 1x1 depth reductions, depth-wise concatenations
 - better results with shallower networks

https://towardsdatascience.com/the-w3h-of-alexnet-vggnet-resnet-and-inception-7baaaecccc96

https://arxiv.org/abs/1512.03385 (He et al., 2016) (Szegedy et al., 2015) http://arxiv.org/abs/1409.4842



weight laver

weight layer

relu

↓ relu

 $\mathcal{G}(\mathbf{x})$

projection

identity

 $\mathcal{F}(\mathbf{x})$

 $\mathcal{F}(\mathbf{x}) + \mathcal{G}(\mathbf{x})$



Pretrained CNNs

https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/

• Faster R-CNN

- object detection harder task
- detecting boxes (regions) for multiple objects in image
- Pipeline:
 - Region prediction network (detect salient boxes)
 - Region-of-interest pooling (consolidate features)
 - Region-based CNN (classify)



Region prediction

- pretrained VGG as feature extraction
 - features for each of the anchor points (regularly spaced in the image)
- for each anchor point, predict:
 - anchor base size & h/w ratio (e.g. 64-128-256px, 0.5/1/1.5)
 - p(this is object) & p(this is background)
 - anchor Δx , Δy , Δh , Δw
 - all of this via convolutions 😳
- trained using object/non-object anchors
- overlapping predictions unified

(Ren et al., 2015) https://arxiv.org/abs/1506.01497

https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/







R-CNN classification

NPFL099 L11 2022

- basically the same as image classification (given region)
 - with one more box coordinates fix
- sharing VGG features from RPN
 - this makes it much faster (only the pooling & prediction layers are new)



https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/

Using Transformers

(Carion et al., 2020) https://arxiv.org/abs/2005.12872

- Object detection: CNN + Transformer
 - trained end-to-end
 - predicts all objects at once
 - max. *k* objects: class + box (or "null")
 - set-based loss (any order allowed)
- Classification: Transformer-only
 - split image into 16x16 patches
 - flatten (256xRGB=768-dim)
 - pass through a single "embedding" matrix
 - feed into standard transformer
 - adding positional embeddings
 - special "start" token
 - attention here used in a feed-forward classifier

(Dosovitskiy et al., 2020)

http://arxiv.org/abs/2010.11929 https://youtu.be/TrdevFK_am4





Visual Dialogue

- **Task**: have a meaningful dialogue about an image
 - close to visual QA: human asks, system responds
 - but VD is multi-turn & human doesn't see the image (just a caption)
 - follow-up questions possible coreference
 - people are not primed by the image when asking questions
- not much realistic purpose other than to test the models
 - dataset of 10-turn dialogues on 120k images
 - collected via crowdsourcing
 - connecting 2 people live to talk about an image
 - not very deep dialogue: history only needed in ~11% cases

(Agarwal et al., 2020) https://www.aclweb.org/anthology/2020.acl-main.728



Caption: A man and woman on bicycles are looking at a map. Person A (1): where are they located Person B (1): in city Person A (2): are they on road Person B (2): sidewalk next to Person A (3): any vehicles Person B (3): 1 in background Person A (4): any other people Person B (4): no Person A (5): what color bikes Person B (5): 1 silver and 1 vellow Person A (6): do they look old or new Person B (6): new bikes Person A (7): any buildings Person B (7): yes Person A (8): what color Person B (8): brick Person A (9): are they tall or short Person B (9): i can't see enough of them to tell Person A (10): do they look like couple

Person B (10): they are



A4: nobody is in the room Q5: can you see on the outside ? A5: no, it is only inside Q6: what color is the sink ? A6: the sink is white Q7: is the room clean ? A7: it is very clean Q8: is the toilet facing the sink ? A8: yes the toilet is facing the sink Q9: can you see a door ? A9: yes, I can see the door Q10 what color is the door ?

Caption: A10 the door is tan colored A sink and toilet in a small room.

NPFL099 L11 2022

Base Visual Dialogue Models

(Das et al., 2017) http://arxiv.org/abs/1611.08669



Visual Dialogue Evaluation

- BLEU etc. possible but not used here
- IR setup used instead
 - system given ground-truth dialogue history + user input & 100 candidate answers to score/rank
- IR metrics:
 - ground-truth response rank (average)
 - **recall@k** (% cases where ground-truth is included in top k)
 - mean reciprocal rank: $\oint \frac{1}{\text{ground truth rank}} (1 \text{ if ground truth is first, } 0.5 \text{ if second etc.})$
 - normalized discounted cumulative gain
 - for multiple acceptable answers out of the 100 candidates
 - DCG: $\sum_{i=1}^{100} \frac{c_i \text{ relevant?}}{\log_2(i+1)}$, normalize by highest possible DCG (all good answers on top)
- problem: images only give modest gain over text-only models

<u>https://en.wikipedia.org/wiki/Discounted_cumulative_gain</u>

https://visualdialog.org/challenge/2019#evaluation

https://medium.com/@_init_/notes-on-the-ndcg-metric-used-in-the-visual-dialog-challenge-2019-90cf443b93dc

Guess What

(Strub et al., 2017) https://www.ijcai.org/proceedings/2017/385

- guessing one of the objects in an image
 - GuessWhat data (150k guessing dialogues)
- 3 models:

NPFL099 L11 2022

- question generation LSTM
 - running through all previous questions
 - conditioned on VGG image features & previous replies
- "oracle" reply generation (Y/N/NA)
 - feed-forward from LSTM question encoding
 + object category + object size & position in image
- guesser select object from list of candidates-
 - dot product & softmax over last LSTM generator state + candidate objects categories & sizes/positions
 - triggered at the end of the dialogue
- trained jointly with RL

Is it a person?

Is it a snowboard?

Is it the red one?

person in blue?

Is it an item being worn or held?

Is it the one being held by the

Yes

Yes

No

Yes



Is it a cow?	Ye
Is it the big cow in the middle?	No
Is the cow on the left?	No
On the right ?	Yes
First cow near us?	Ye

Shopping Dialogues

(Saha et al., 2018) http://arxiv.org/abs/1704.00200

- a more practical task: fashion e-shop
 - both user & system can show images
 - user can reference/talk about images
- semi-automatic data (150k dialogues)
 - 84 predefined intents
 - generated dialogue flows
 - based on predefined state automata over intents
 - utterances for each intent written by annotators
 - backed by a large manually curated DB of fashion items with images (1M items)
 - taxonomy (man > apparel > lower body > trousers)
 - (generated) celebrity profiles/endorsements

SHOPPER: Hello

AGENT: Hello, is there anything i can help you with today? SHOPPER: show me some espadrilles in beige/blue within \$200. AGENT: Sure. let me just quickly browse through my catalogue



AGENT: What you can see in the fifth item, michael michael kors brings you an effortless weekend look with the wrapped espadrille heel and crisscross styling of these dana flats. And what you see in the 3rd item, steve madden, the company, was founded by steve madden, the man, in 1990 in brooklyn, where he started out making trend-advancing shoes one pair at a time.

SHOPPER: List more in the upper material as in the 5th image and style as the 3rd and the 5th



AGENT: The 2nd one, viscata espadrilles are 100% authentic and handcrafted in spain inspired by the mediterranean, viscata shoes are lightweight and made of natural jute rope, premium leathers, suedes and woven cloth.

SHOPPER: Can you list images of the 1st result from all orientations?



AGENT: Image from the front, right and back orientations respectively

SHOPPER: And the 3rd one?



AGENT: Image from the front, right, back and left orientations respectively

Id	Automata	State Description
	State Type	
1	greeting	Shopper greets
2	self-info	Shopper gives information about him/herself
3	give-criteria	describes his requirements
4	show-image	Agent shows relevant responses for the Shopper's query
5	give-image-	Agent generates short description of the product, using
	description	visual and catalog information
6	Like/Dislike	Shopper expresses negative or positive preference specif-
	specific	ically towards one or more items previously or currently
	items / over-	shown, or a overall general preference towards all the
	all items,	items and optionally shows a new image to possibly
	show-more	modify his requirements and wants to see more
7	show-	Shopper wants to see an item from different orientations
	orientation	
8	show-	Shopper wants to see similar to a particular item
	similar	

Shopping Dialogues

- Models similar to visual dialogue
 - variants of multimodal HRED
 - VGG image input
- image input
 - turn-level
 - concatenated with utterance
 - seems to work better (fewer turns)
- text/image responses
 - shared encoder
 - text generation (word-by-word)
 - image ranking (needs rough retrieval)
 - so far just "select 1 out of 5"



Using Images to boost NLU

- Grounding all words to images "vokenization"
 - images + captions, nearest neighbor search
 → assigning an image to each word voken
 - train a vokenizer (LM assigning images to tokens) on this
 - apply it to vokenize large training data for BERT
- Finetuning BERT with:
 - masked language modelling (as usual)
 - predicting a voken for all words (masked or not)
 - classification mimicking the vokenizer
- Further finetuning for a NLU task
 - better performance when vokens were used





(Tan and Bansal, 2020)

Situated Tasks (Amazon SimBot Challenge)

- Home tasks commander & follower
 - find best trajectory of actions to complete task
 - include dialogue, ask if unsure
- EMMA integrated architecture for this
 - core: Transformer LM (+sparse attention)
 - "visual tokens" to represent video input
 - explicit actions / words on the output
 - pretraining on many vision & language tasks
 - captioning, image QA, location, relations



(Padmakumar et al., 2022) https://arxiv.org/abs/2110.00534

(Suglia et al., 2022) https://aclanthology.org/2022.sigdial-1.62/



Domain Adaptation

- pretraining
 - BERT, but also any other model
 - weight sharing: copy weights for similar slots in target domain
- delexicalization
 - assuming your domains are similar (e.g. TVs → PCs)
- pseudo in-domain data selection
 - find data similar to your domain in the source domain
- forcing shared latent space (see few-shot end-to-end models)
- multi-task training
 - your task in source domain & different task in target domain
- partial handcrafting (see Hybrid Code Networks)

Summary

- "traditional" multimodal systems, with components
 - combination of off-the-shelf components
 - parallels for ASR/NLU & NLG/TTS in I/O modalities
 - dialogue typically quite simple
 - modalities: static graphics / touch / gaze / facial expr. / avatars / robots
 - often support multi-party dialogue
- end-to-end multimodal systems
 - mostly experimental, based on HRED with pretrained CNNs
 - VGG, ResNet, Inception (just image classification), Faster R-CNN (+object detection)
 - visual dialogue: questions & answers about an image
 - guessing: finding an object in image
 - task-oriented: shopping dialogue with product images
 - situated tasks: discussing & executing household actions

Thanks

Contact us:

<u>https://ufaldsg.slack.com/</u> {odusek,hudecek,kasner}@ufal.mff.cuni.cz Skype/Meet/Zoom (by agreement)

Get the slides here:

http://ufal.cz/npfl099

References/Inspiration/Further:

- Volha Pethukova's course (Uni Saarland): <u>https://www.lsv.uni-saarland.de/multimodal-dialogue-systems-summer-2019/</u>
- McTear et al. (2016): The Conversational Interface Talking to Smart Devices
- Delgado & Araki (2005): Spoken, Multilingual and Multimodal Dialogue Systems: Development and Assessment
- papers referenced on slides

Next week: last lecture & labs

Exam

- In-person written test, 10 questions covering lectures, 10 points each
 - 50% on homework assignments needed to do the test
 - counts for 75% of the grade, 25% comes from homework assignments
 - grades: 1 = 87%+, 2 = 74%+, 3 = 60%+ (for the weighted combo)
 - expected 1 hr, but you'll be given at least 2hrs (no pressure on time)
- Question type: 2-3 sentences to answer
 - explanation of terms/concepts
 - no exact formulas needed (if needed, they might be provided)
 - but you should know the principles of how stuff works
 - relationships between concepts ("what's the difference between X & Y")
 - "how would you build X"
 - focused on "important" stuff see summaries at the end of each lecture
 - list of possible questions to be published soon (by Dec 31)