NPFL099 Statistical Dialogue Systems 8. Natural Language Generation

http://ufal.cz/npfl099

Zdeněk Kasner, Ondřej Dušek & Vojtěch Hudeček 21.11.2022



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



Natural Language Generation

= task of automatically producing text in e.g. English (or any other language)

• covers many subtasks:

	task	input	output	
	unconditional language generation	Ø	arbitrary text	
NLG in a narrow sense	conditional language generation	short text prompt	continuation of the prompt	
	machine translation	text in language A	text in language B	
	summarization	long text	text summary	
	question answering	question	answer	
	image captioning	image	image caption	
	data-to-text generation	structured data	description of the data	
	dialogue response generation	dialogue act	system response	

NLG Objectives

• general NLG objective:

given input & communication goal create accurate + natural, well-formed, human-like text

• additional NLG desired properties:

- variation (avoiding repetitiveness)
- simplicity (saying only what is intended)
- adaptability (conditioning on e.g. user model)

NLG in Dialogue Systems

• in the context of dialogue systems:

NLG: system action → system response

"what the system wants to say"

"actually saying it"

- system action
 - selected by the dialogue manager
 - may be conditioned on:
 - dialogue state
 - dialogue history (→ referring expressions, avoiding repetition)
 - user model (→ "user wants short answers")

NLG Subtasks (Textbook Pipeline) = how proper NLG had to be done before neural approaches



NLG Subtasks (Textbook Pipeline)

Example: classical NLG pipeline in medical domain



NLG Basic Approaches

hand-written prompts ("canned text")

- most trivial hard-coded, no variation
- doesn't scale (good for DTMF phone systems)

• templates ("fill in blanks")

- simple, but much more expressive covers most common domains nicely
- can scale if done right, still laborious
- most production dialogue systems

• grammars & rules

- grammars: mostly older research systems
- rules: mostly content & sentence planning

machine learning

- modern research systems
- pre-neural attempts often combined with rules/grammar
- NNs made it work much better









Template-based NLG

- most common in commercial dialogue systems
- simple, straightforward, reliable
 - custom-tailored for the domain
 - complete control of the generated content
- lacks generality and variation
 - difficult to maintain, expensive to scale up
- can be enhanced with rules
 - e.g. articles, inflection of the filled-in phrases
 - template coverage/selection rules (heuristics, random variation)
- can be a good starting point for ML algorithms
 - post-editing / reranking the templates with neural language models



Template-based NLG – Examples

Example: Facebook

{user} shared {object-owner}'s {=album} {title}
Notify user of a close friend sharing content

* {user} is female. {object-owner} is not a person or has an unknown gender.
{user} sdilela {=album} "{title}" uživatele {object-owner}
{user} sdilela {object-owner} uživatele {=album}{title}
* New translation

(Facebook, 2015)



(Facebook, 2019)

Template-based NLG – Examples

Example: Dialogue assistants

Alexa

On the **Intents** detail page, the **Intent Slots** section after the **Sample Utterances** section displays the slots you add. When you highlight a word or phrase in an utterance, you can add a new slot or select an existing slot.

For example, the set of utterances shown earlier now looks like the following example.

i am going on a trip on {travelDate}
i want to visit {toCity}
I want to travel from {fromCity} to {toCity} {travelDate}
I'm {travelMode} from {fromCity} to {toCity}
i'm {travelMode} to {toCity} to go {activity}

(https://developer.amazon.com/en-US/docs/alexa/custom-skills/ create-intents-utterances-and-slots.html)

Mycroft

Order some {food}.
Order some {food} from {place}.
Grab some {food}.
Grab some {food} from {place}.

Rather than writing out all combinations of possibilities, you can embed them into a single line by writing each possible option inside parentheses with | in between each part. For example, that same intent above could be written as:

(Order | Grab) some {food} (from {place} |)

(https://mycroft-ai.gitbook.io/docs/mycroft-technologies/padatious)

Template-based NLG – Examples

Example: Research systems

'iconfirm(to_stop={to_stop})&iconfirm(from_stop={from_stop})':
 "Alright, from {from_stop} to {to_stop},",

'iconfirm(to_stop={to_stop})&iconfirm(arrival_time_rel="{arrival_time_rel}")':
 "Alright, to {to_stop} in {arrival_time_rel},",

'iconfirm(arrival_time="{arrival_time}")':
 "You want to be there at {arrival_time},",

'iconfirm(arrival_time_rel="{arrival_time_rel}")':
 "You want to get there in {arrival_time_rel},",

(Alex public transport information rules) https://github.com/UFAL-DSG/alex

CONFIRM!!date!!@	The date is @.		
CONFIRM!!party_size!!@	The reservation is for @ people.		
CONFIRM!!restaurant_name!!@	Booking a table at @.		
CONFIRM!!time!!@	The reservation is at @.		
GOODBYE	Have a good day.		
INFORM!!cuisine!!@	They serve @ kind of food.		
INFORM!!has_live_music!!False	They do not have live music.		
INFORM!!has_live_music!!True	They have live music.		



Is it a good restaurant? Are they expensive?

User

(Kale & Rastogi, 2020) https://www.aclweb.org/anthology/2020.emnlp-main.527

Grammar / Rule-based NLG

- based on top of linguistic theories
- state-of-the-art research systems until NLG the arrival of NNs
- rules for building tree-like structures
 → rules for tree linearization
- reliable, more natural than templates
- takes a lot of effort, naturalness still not human-level
- see NPFL123 for more details





Neural NLG

- learning the task from the data
- sequence-to-sequence generation / editing / re-ranking
- fluency can match human-level, minimal hand-crafting
- not controllable ("black-box"), semantic inaccuracies (omissions / hallucinations), low diversity, expensive data gathering, expensive training, expensive deployment

 \rightarrow promising research area 😉



encoder-decoder

- *RNN:* encoder updates the hidden state → decoder is initialized with the hidden state
- *Transformer*: encoder generates a sequence of hidden states → decoder attends to this sequence
- pretrained Transformers: BART, T5 (trained on sequence denoising)

decoder-only

- input sequence is prepended as a context, the decoder generates continuation
- pretrained Transformer models: GPT-2, GPT-3 (trained on autoregressive language modelling)
- training vs. inference:

Encoder-Decoder Training



Encoder-Decoder Inference



- for each time step t, the decoder outputs a probability distribution: $P(y_t | y_{1:t-1} \mathbf{X})$
- how to use it?
- exact inference: find a sequence maximizing $P(y_{1:T} | \mathbf{X})$
 - not possible in practice (why? and is it our goal?)
- approximation algorithms
 - greedy search
 - beam search
- stochastic algorithms
 - random sampling
 - top-k sampling
 - nucleus sampling (=top-p sampling)

(+ repetition penalty → decreasing probabilities of generated tokens)

- greedy search: always take the argmax
 - does not necessarily produce the most probable sequence (why?)
 - often produces dull responses



Example:

Context: Optimal Response : Greedy search: Try this cake. I baked it myself. This cake tastes great. This is okay.

many examples start with "This is", no possibility to backtrack

- **beam search:** try *k* continuations of *k* hypotheses, keep *k* best
 - better approximation of the most probable sequence, bounded memory & time
 - allows re-ranking generated outputs
 - \circ k=1 \rightarrow greedy search



Reranking:



(Ondřej's PhD thesis, Fig. 7.7) http://ufal.mff.cuni.cz/~odusek/2017/docs/thesis.print.pdf

- **top-k sampling:** choose top k options (~5-500), sample from them
 - \circ avoids the long tail of the distribution
 - more diverse outputs



- **nucleus sampling:** choose top options that cover >= *p* probability mass (~0.9)
 - \circ "k" is adapted according to the distribution shape



RNN-based Approaches: TGen

RNN seq gen + classif

- standard LSTM with attention
 - encoder triples <intent, slot, value>
 - decodes words (possibly delexicalized)



• beam search & reranking

- DA classification of outputs
- checking against input DA



RNN-based Approaches: RNNLM

- using enhanced LSTM cells (SC-LSTM)
 - special gate to control slot mentions
- autoregressive generation
- conditioned on DA represented as binary vector
 - generating delexicalized texts
- domain adaptation
 - replacing delexicalized slots
 - very related domains only

(Wen et al, 2015; 2016) http://aclweb.org/anthology/D15-1199 http://arxiv.org/abs/1603.01232



Delexicalization Alternatives

• **copy mechanism** (see NLU & the next slide)

- select (or interpolate) between:
 - generating a new token
 - copying a token from input
- does away with the pre/postprocessing

inflection model

- useful for languages with rich morphology (e.g., Czech)
- a simple LM such as RNN LM

pretrained models

- the model learns to copy and inflect words during pretraining → may work implicitly on downstream tasks
- works well for high-resource languages

inform(name=Baráčnická rychta, area=Malá Strana)



(Dušek & Jurčíček, 2019) https://arxiv.org/abs/1910.05298

Delexicalization Alternatives – Copy Mechanism



Pretrained LMs

- GPT-2 or BART / T5 finetuned for NLG
 - different pretraining tasks similar outcomes
- works nicely when simply finetuned for data-to-text
 - encode linearized data, decode text
 - the model learns copying implicitly
- mBART / mT5 (multilingual) → allows multilingual generation
 - can generate e.g. Russian outputs from English triples
- are we done now?



(Kale & Rastogi, 2020) https://www.aclweb.org/anthology/2020.in lg-1.14

(Liu et al., 2020) http://arxiv.org/abs/2001.08210

(Kasner & Dušek, 2020)

webnlg-1.20/

https://aclanthology.org/2020.

NPFL099 L8 2022

Pretrained LMs with Reranking

- aiming to improve semantic accuracy
- seq2seq + reranking with GPT-2 & RoBERTa
- GPT-2 fine-tuned for <data> name[Zizzi] eatType[bar] <text> Zizzi is a bar.

prompt (fed into GPT-2) decoded given the prompt

- beam search decoding
- RoBERTa for classification
 - accurate/omission/repetition/hallucination/value error
 - training data synthesized
 - "accurate" examples from original training data
 - others created by manipulating the data and texts (adding/removing/replacing sentences and/or data items)



Few-shot NLG with Pretrained LMs

- learning from very few (less than ~200) training examples
- GPT-2 with a copy mechanism
 - LM fine-tuned, forced to copy inputs
 - additional loss term for copying
- retrieving "prototypes" guiding the generator
 - prototype: most similar exemplar according to BERT cosine similarity
 - prototype concatenated with the input
- few-shot prompting
 - prepending a few (~3) input-output examples as a context
 - generating the output with GPT-2 XL
 - no finetuning

(Keymanesh et al., 2022) http://arxiv.org/abs/2205.11505 (Chen et al., 2020) https://www.aclweb.org/anthology/2 020.acl-main.18/

(Su et al., 2021) http://arxiv.org/abs/2108.12516

Few-shot Prompt
Translate Graph to English:
Graph: <i><h></h></i> Paulo Sousa <i><r></r></i> CLUB <i><t></t></i> ACF Fiorentina
English: Paulo Sousa plays for ACF Fiorentina. ####
Graph: <i><h></h></i> Dave Challinor <i><r></r></i> CLUB <i><t></t></i> Colwyn Bay F.C.
English: Dave Challinor plays for Colwyn Bay F.C. ####
Graph: <i><h></h></i> Alan Martin (footballer) <i><r></r></i> CLUB <i><t></t></i> Hamilton Academical F.C.
English:

NPFL099 L8 2022

Few-shot NLG: Templates + Pretrained LM

- start with simple templates (one piece per system action)
 - a bit of handcrafting, but manageable for many datasets
- concatenate them and then use pretrained LMs (e.g. T5/BART) to rephrase them
 - basically text-to-text denoising, i.e. what the models were originally trained to do
- needs less data & generalizes to new domains

templates

CONFIRM!!date!!@	The date is @.		
CONFIRM!!party_size!!@	The reservation is for @ people.		
CONFIRM!!restaurant_name!!@	Booking a table at @.		
CONFIRM!!time!!@	The reservation is at @.		
GOODBYE	Have a good day.		
INFORM!!cuisine!!@	They serve @ kind of food.		
INFORM!!has_live_music!!False	They do not have live music.		
INFORM!!has_live_music!!True	They have live music.		

system



(Kale & Rastogi, 2020) <u>https://www.aclweb.org/anthology</u> /2020.emnlp-main.527

Zero-shot NLG: Templates + Pretrained LM

NLG without human-written references

- start with templates → postprocess them with BART-based models
 - trained for text-based operations learned from Wikipedia



• improvement: using prompted GPT-3 instead of hand-crafted templates

(Xiang et al., 2022) http://arxiv.org/abs/2210.04325

Content Planning: Content Selection

explicit content selection

- usually done by DM in dialogue systems
- needed for complex inputs, e.g. sports report generation
 - records (team / entity / type / value) → summary
 - content selection: pointer network
- still largely unsolved problem w.r.t. semantic accuracy

(Thomson & Reiter, 2022) http://arxiv.org/abs/2108.05644



http://arxiv.org/abs/1809.00582

Text Generation p(y|r,z)Plan SOS Va V1 Attention **Content Plan** Encoding r p(z|r) r_2^{cs} $r^{cs}_{|r|}$ Content r_4^{cs} Selection & Planning $r_{|r|}$ EOS Decoder Vector 🔿 Content Selection Gate Encoder

 $p_{gen}(y_4|r, z, y_{<4})$ $p_{copy}(y_4|r, z, y_{<4})$

datt

content plan

(for the 1st sentence)

Example of NLG with content planning

source statistics (excerpt)

TEAM	WIN	LOSS	PTS	FG	PCT	RB	AST	14	3
Pacers	4	6	99	8	42	40	17		
Celtics	5	4	105	12	44	47	22	92	
PLAYE	R	H/V	AST	RB	PTS	FG	CITY		
Jeff Teague		н	4	3	20	4	India	na	
Miles Ti	Ailes Turner		1	8	17	6	India	na	
saiah Thomas		s V	5	0	23	4	Bosto	n	
Kelly O	lynyk	V	4	6	16	6	Bosto	n	141
Amir Johnson		V	3	9	14	4	Bosto	n	

PTS: points, FT_PCT: free throw percentage, RB: rebounds, AST: assists, H/V: home or visiting, FG: field goals, CITY: player team city.

Value	Entity	Туре	H/V	
Boston	Celtics	TEAM-CITY	V	
Celtics	Celtics	TEAM-NAME	V	
105	Celtics	TEAM-PTS	V	
Indiana	Pacers	TEAM-CITY	н	
Pacers	Pacers	TEAM-NAME	н	
99	Pacers	TEAM-PTS	н	
42	Pacers	TEAM-FG_PCT	н	
22	Pacers	TEAM-FG3_PCT	н	
5	Celtics	TEAM-WIN	V	
4	Celtics	TEAM-LOSS	V	
Isaiah	Isaiah_Thomas	FIRST_NAME	V	
Thomas	Isaiah_Thomas	SECOND_NAME	V	
23	Isaiah_Thomas	PTS	V	
5	Isaiah_Thomas	AST	V	
4	Isaiah_Thomas	FGM	V	
13	Isaiah_Thomas	FGA	V	
Kelly	Kelly_Olynyk	FIRST_NAME	V	
Olynyk	Kelly_Olynyk	SECOND_NAME	V	
16	Kelly_Olynyk	PTS	V	
6	Kelly_Olynyk	REB	V	
4	Kelly_Olynyk	AST	V	

target text

The **Boston Celtics** defeated the host **Indiana Pacers 105-99** at Bankers Life Fieldhouse on Saturday. In a battle between two injury-riddled teams, the Celtics were able to prevail with a much needed road victory. The key was shooting and defense, as the **Celtics** outshot the **Pacers** from the field, from three-point range and from the free-throw line. Boston also held Indiana to **42 percent** from the field and **22 percent** from long distance. The Celtics also won the rebounding and assisting differentials, while tying the Pacers in turnovers. There were 10 ties and 10 lead changes, as this game went down to the final seconds. Boston (**5–4**) has had to deal with a gluttony of injuries, but they had the fortunate task of playing a team just as injured here. **Isaiah** Thomas led the team in scoring, totaling **23 points and five assists on 4–of–13** shooting. He got most of those points by going 14–of–15 from the free-throw line. **Kelly Olynyk** got a rare start and finished second on the team with his **16 points, six rebounds and four assists**.

> (Puduppully et al., 2019) http://arxiv.org/abs/1809.00582

Content Planning: Ordering & Aggregation

ordering the facts + aggregating them into sentences

- content already selected at this point
- can help the generator not to miss any facts
- for graphs with oriented edges:
 - generating all possible content plans using DFS (possibly pruning unpromising branches) → re-ranking the plans using a feature-based classifier
- for a set of key-value pairs:
 - using Conditional Random Field (CRF) for finding the optimal plan



(Moryossef et al., 2019a,b) http://arxiv.org/abs/1904.03396 https://arxiv.org/pdf/1909.09986.pdf



Realizing from Trees

seq2seq + copy seq gen

• NLG with tree-shaped inputs

- simple case: discourse relations (discourse connectives, sentence splits) between individual fields
 - much flatter than usual syntactic trees
- improvements to account for the input structure:
 - constrained beam search decoding, tree-LSTM, self-training on synthetic data

Reference 1 Reference 2 E2E MR	JJ's Pub is not family friendly, but has a high customer rating of 5 out of 5. It is a restaurant near the Crowne Plaza Hotel. JJ's Pub is not a family friendly restau- rant. It has a high customer rating of 5 out of 5. You can find it near the Crowne Plaza Hotel. name[JJ's Pub] rating[5 out of 5] familyFriendly[no] eatType[restaurant] near[Crowne Plaza Hotel] CONTRAST [INFORM [name[JJ's Pub]	Query When will it snow next?	Context Reference date: 29th September 2018	MR [CONTRAST [INFORM_1] [LOCATION [CITY Parker]] [CONDITION_NOT snow] [DATE_TIME [DAY 29] [MONTH September] [YEAR 2018]]] [INFORM_2 [DATE_TIME [DAY 29] [MONTH September] [YEAR 2018]] [LOCATION [CITY Parker]] [CONDITION heavy rain showers] [CLOUD_COVERAGE partly cloudy]]	Response Parker is not expecting any snow, but today there's a very likely chance of heavy rain showers, and it'll be partly cloudy	
Our MR for Reference 1	familyFriendly[no]] INFORM [rating[5 out of 5]]] INFORM [eatType[restaurant] near[Crowne Plaza Hotel]]	Annotated Response [CONTRAST [INFORM_1 [LOCATION [CITY Parker]] is not expecting any [CONDITION_NOT snow]], but [IN- FORM_2 [DATE_TIME [COLLOQUIAL today]] there's a [PRECIP_CHANCE_SUMMARY very likely chance] of [CONDITION heavy rain showers] and it'll be [CLOUD_COVERAGE partly cloudy]]				

Data Noise & Cleaning

- NLG errors are often caused by **data errors**
 - ungrounded facts (← hallucinating)
 - missing facts (← forgetting)
 - domain mismatch
 - noise (e.g. source instead of target)
 - just 5% untranslated stuff kills an NMT system
- easy-to-get data are noisy
 - web scraping lot of noise, typically not fit for purpose
 - crowdsourcing workers forget/don't care
- cleaning improves situation a lot
 - can be done semi-automatically up to a point

(Khayrallah & Koehn, 2018) https://www.aclweb.org/anthology/W18-2709

Original MR and an accurate reference

 $\label{eq:mass_share} \begin{array}{l} \textbf{MR} & \texttt{name}[\texttt{Cotto}], \ \texttt{eatType}[\texttt{coffee shop}], \ \texttt{food}[\texttt{English}], \ \texttt{priceRange}[\texttt{less} \\ \texttt{than} \ \pounds20], \ \texttt{customer_rating}[\texttt{low}], \ \texttt{area}[\texttt{riverside}], \ \texttt{near}[\texttt{The Portland Arms}] \end{array}$

Reference At the riverside near The Portland Arms, Cotto is a coffee shop that serves English food at less than $\pounds 20$ and has low customer rating.

Example corrections

Reference: Cotto is a coffee shop that serves English food in the city centre. They are located near the Portland Arms and are low rated. Correction: removed price range; changed area Reference: Cotto is a cheap coffee shop with one-star located near The Portland Arms. Correction: removed area

A faulty correction

Reference: Located near The Portland Arms in riverside, the Cotto coffee shop serves English food with *a price range of \$20* and a low customer rating.

Correction: incorrectly(!) removed price range – our script's slot patterns are not perfect

(Dušek et al., 2019) https://arxiv.org/abs/1911.03905

(Wang, 2019) https://www.aclweb.org/anthology/W19-8639/

Data Augmentation

using synthetic data for improving model performance and robustness

- quite hard for NLG
- where to get the data:
 - extracting (noisy) structured input from unlabeled text
 - keywords, information extraction
 - recombination of existing data inputs
 - model-generated outputs → filtering based on cycle consistency
 - paraphrasing existing data outputs
- how to apply it:
 - task-specific pretraining on synthetic data
 - mixing synthetic data with the training data

(Elder et al., 2020) https://www.aclweb.org/anthology/2020.a cl-main.665 (Montella et al., 2020) https://arxiv.org/pdf/2012.00571.pdf (Chang et al., 2021) http://arxiv.org/abs/2102.03556 (Lee et al. 2021) http://arxiv.org/abs/2110.06800



Summary

- NLG: system action → system response
- templates work pretty well
- **seq. generation** with pretrained Transformer LMs = best among data-driven
- problems hallucination, not enough diversity, needs lots of data
 - reranking (with NLU), delexicalization/copy nets, pretrained models
 - different decoding styles (sampling)
 - NLG-NLU joint training
 - for data cleaning, self-training, semi-supervised training
 - 2-step: planning & realization
 - tree decoding more supervision
 - few-shot with pretrained LMs
 - data manipulation: cleaning, augmentation

Thanks

Contact us:

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

Get these slides here:

http://ufal.cz/npfl099

References/Inspiration/Further:

- Gatt & Krahmer (2017): Survey of the State of the Art in Natural Language Generation: Core tasks, applications and evaluation <u>http://arxiv.org/abs/1703.09902</u>
- Sharma et al. (2022). Innovations in Neural Data-to-text Generation. <u>https://arxiv.org/pdf/2207.12571.pdf</u>
- Ondřej's PhD thesis (2017), especially Chapter 2: <u>http://ufal.mff.cuni.cz/~odusek/2017/docs/thesis.print.pdf</u>

Icons from https://www.flaticon.com/

Labs in 10 minutes Assignment 4

Next week: End-to-end models