

Internship: master or engineering degree

Deep Learning for Data to Text and Text to Data Generation:

Adapting large-scale Transformer models to language generation from data and data generation from textual data.

Contact: patrick.gallinari@sorbonne-universite.fr, t.ricatte@criteo.com

When: starting from March to May 2021 for 5 to 6 months

Location: SCAI, Sorbonne Center for Artificial Intelligence, Sorbonne University, Paris

Gratification: classical internship gratification around 550 E/ month

Skills: computer science or applied math profile with a strong background in machine learning and/or natural language processing.

Context

Text generation has witnessed impressive progress with recent developments of Deep Learning (DL) methods. DL language models initially based on Recurrent Neural Networks and more recently on transformer architectures can be trained on huge quantities of data and then be used for generating high quality text, conditioned on some initial input. These systems currently achieve impressing results in different text generating tasks including free text generation (Radford et al. 2018, Raffel et al. 2020) , abstractive summarization (Scialom et al 2020), or semi structured text generation (Bien et al. 2020).

Data to text and text to data

Knowledge sources are often encoded into structured format such as indexes, tables, triplets, ontologies, knowledge bases, or even raw numerical data. These data are easily readable by machines, but hardly interpretable by humans. On the opposite, textual information, easily accessible to humans is often complex to exploit by machines. A key challenge and an emerging field in machine learning and natural language processing, is the transcription of structured data to text and the inverse problem of transforming raw text into structured data.

The former problem is called data-to text generation and it occurs in several applications like journalism, medical diagnosis, financial reports, sport broadcasts, dialogue oriented tasks - generating responses, summarization. It may also be a component of explainable AI systems where by allowing humans to interact with machines. Data can come in different formats such as tables, graphs, etc. The RotoWire task for example consists in generating NBA game summaries from tables providing the game statistics (Figure 1).

TEAM	H/V	WINS	LOSSES	PTS	REB	AST	...
Hawks	H	46	12	95	42	27	...
Magic	V	19	41	88	40	22	...

PLAYER	PTS	REB	AST	STL	BLK	CITY	...
Al Horford	17	13	4	2	0	Atlanta	...
Kyle Korver	8	3	2	1	2	Atlanta	...
Jeff Teague	17	0	7	2	0	Atlanta	...
N. Vučević	21	15	3	1	1	Orlando	...
Tobias Harris	15	4	1	2	1	Orlando	...
...

H/V: home or visiting; PTS: points; REB: rebounds;
AST: assists; STL: steals; BLK: blocks

The **Atlanta Hawks (46-12)** beat the **Orlando Magic (19-41)** **95-88** on Friday. **Al Horford** had a good all-around game, putting up **17 points, 13 rebounds, four assists and two steals** in a tough matchup against **Nikola Vučević**. **Kyle Korver** was the lone Atlanta starter not to reach double figures in points. **Jeff Teague** bounced back from an illness, he scored **17 points** to go along with **seven assists and two steals**. After a rough start to the month, the **Hawks** have won three straight and sit atop the Eastern Conference with a nine game lead on the second place Toronto Raptors. The **Magic** lost in devastating fashion to the Miami Heat in overtime Wednesday. They blew a seven point lead with 43 seconds remaining and they might have carried that with them into Friday's contest against the **Hawks**. **Vučević** led the **Magic** with **21 points and 15 rebounds**. **Aaron Gordon** (ankle) and **Evan Fournier** (hip) were unable to play due to injury. The **Magic** have four teams between them and the eighth and final playoff spot in the Eastern Conference. The **Magic** will host the Charlotte Hornets on Sunday, and the **Hawks** with take on the Heat in Miami on Saturday.

Figure 1: example from the Rotowire Corpus (Wiseman et al. 2017). Given the game statistics (left), the objective is to generate a textual summary (right).

As a second example, the WebNLG challenge (Ferreira et al. 2020) consists in transcribing graphs composed of RDF triplets to text and vice versa (Figure 2).

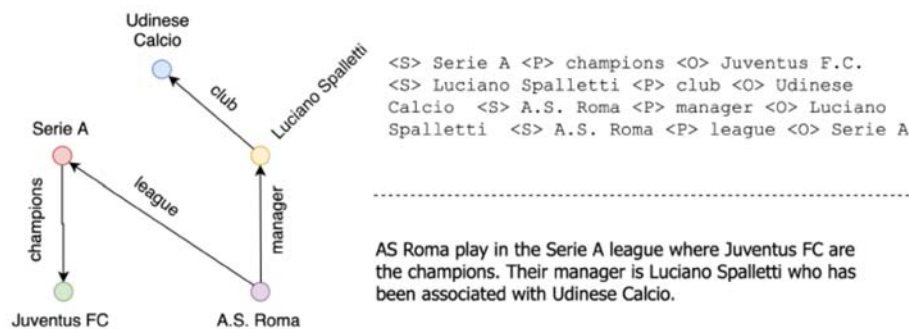


Figure 2: example from the WebNLG corpus (figure from (Kale 2020)). On the left a graph structure, on the right top its serialization as a set of RDF triples, on the right bottom, a possible associated text.

The latter problem is known as *semantic parsing* and comes in different instantiations like information retrieval, reasoning over the structured data (table or graph), generating symbolic queries (e.g. SQL) from text which may be used for example in dialog systems when the answer to a user requires querying a database, or generating abstract meaning representations (Belavicqua 2021). The 2020 edition of the WebNLG challenge (Figure 2) addressed the two tasks of data to text and semantic parsing.

Internship objective

Recent progress has been made for data to text tasks through the use of recent transformers (Agarwal et al. 2020, Kale et al. 2020, Guo et al. 2020) or recurrent sequence to sequence models (Rebuffel et al. 2020). Neural methods trained end to end achieve state of the art performance for different D2T tasks. They however suffer from different pitfalls, like data scarcity (small size annotated corpora), hallucinations (generation of linguistic but non-factual sentences), low coverage of table evidence. They are still restricted to tasks of low complexity. For the dual task of text to data generation, state of the art still relies on complex pipelines integrating different components and sequence-to-sequence models still lag behind the best models. Tasks are also of limited complexity.

The objective of the internship is to adapt recently proposed sequence-to-sequence transformers like Google T5 (Raffel et al. 2020) or Facebook BART (Lewis et al. 2019) for dealing within a unified framework with the dual tasks of text to data (T2D) and data to text (D2T) generation. Initial attempts have been performed in this direction very recently, but are restricted to very specific tasks like Abstract Meaning Representation to Text transcription (Bevilacqua 2021). Besides the adaptation

of existing methods and the development of a new unified framework, a major challenge concerns the scarcity of labeled data. For that, two complementary directions will be considered. One is transfer learning: how general knowledge from models trained in an unsupervised way on large corpora could be adapted to Data2Text and Text2Data. A second one is unsupervised learning: how to use unaligned corpora, which may be more largely available for the two tasks. The focus will be on the analysis of items or product descriptions with an objective of generative textual descriptions from item databases and extracting relevant attributes from textual descriptions.

References

- Agarwal, O., Kale, M., Ge, H., Shakeri, S. and Al-Rfou, R. 2020. Machine Translation Aided Multilingual Data-to-Text Generation and Semantic Parsing. *WebNLG* (2020), 1–5.
- Bevilacqua, M., Blloshmi, R. and Navigli, R. 2021. One SPRING to Rule Them Both : Symmetric AMR Semantic Parsing and Generation without a Complex Pipeline. *AAAI* (2021).
- Bien, M., Gilski, M., Maciejewska, M. and Taisner, W. 2020. RecipeNLG : A Cooking Recipes Dataset for Semi-Structured Text Generation. *INLG* (2020), 22–28.
- Deng, X., Awadallah, A.H., Meek, C., Polozov, O., Sun, H. and Richardson, M. 2020. Structure-Grounded Pretraining for Text-to-SQL. <http://arxiv.org/abs/2010.12773> (2020).
- Devlin, J., Chang, M.W., Lee, K. and Toutanova, K. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL HLT* (2019), 4171–4186.
- Ferreira, T.C., Gardent, C., Ilinykh, N., Lee, C. Van Der, Mille, S., Moussallem, D., Shimorina, A., Loria, C. and Loria, D.L. 2020. The 2020 Bilingual , Bi-Directional WebNLG + Shared Task Overview and Evaluation Results (WebNLG + 2020). (2020).
- Guo, Q., Jin, Z., Wang, Z., Qiu, X., Zhang, W., Zhu, J., Zhang, Z. and Wipf, D. 2020. Fork or Fail: Cycle-Consistent Training with Many-to-One Mappings. <http://arxiv.org/abs/2012.07412> (2020).
- Guo, Q., Jin, Z., Qiu, X., Zhang, W., Wipf, D. and Zhang, Z. 2020. CycleGT: Unsupervised Graph-to-Text and Text-to-Graph Generation via Cycle Training. *WebNLG* (2020).
- Kale, M. 2020. Text-to-Text Pre-Training for Data-to-Text Tasks. *INLG* (2020), 97–102.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V. and Zettlemoyer, L. 2019. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. *ACL* (2019), 7871–7880.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D. and Sutskever, I. 2018. Language Models are Unsupervised Multitask Learners. (2018).
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W. and Liu, P.J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *JMLR*. 21, (2020), 1–67.
- Rebuffel, C., Soulier, L., Scouteeten, G. and Gallinari, P. 2020. A Hierarchical Model for Data-to-Text Generation. *ECIR* (2020), 65–80.
- Scialom, T., Dray, P.-A., Lamprier, S., Piwowarski, B. and Staiano, J. 2020. Discriminative Adversarial Search for Abstractive Summarization. *ICML* (2020).
- Wiseman, S., Shieber, S., Rush, A. 2017: Challenges in data-to-document generation. In: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pp. 2253–2263. Association for Computational Linguistics, Copenhagen, Denmark, September 2017.