# Machine Learning for **Language** and Vision

Seminar SS23
Kickoff Meeting, Apr 18, 2023

Xudong Hong, Ruitao Feng

Saarland University

# GPT-4



GPT-4 | The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.

# The Stats



**ChatGPT Sprints to One Million Users**

Time it took for selected online services to reach one million users

| Service | Launched | Time |
| --- | --- | --- |
| Netflix | 1999 | 3.5 years |
| Kickstarter* | 2009 | 2.5 years |
| Airbnb** | 2008 | 2.5 years |
| Twitter | 2006 | 2 years |
| Foursquare*** | 2009 | 13 months |
| Facebook | 2004 | 10 months |
| Dropbox | 2008 | 7 months |
| Spotify | 2008 | 5 months |
| Instagram*** | 2010 | 2.5 months |
| ChatGPT | 2022 | 5 days |

* one million backers   ** one million nights booked   *** one million downloads
Source: Company announcements via Business Insider/Linkedin

statista

# Artificial General Intelligence?

## Sparks of Artificial General Intelligence:
## Early experiments with GPT-4

Sébastien Bubeck    Varun Chandrasekaran    Ronen Eldan    Johannes Gehrke

Eric Horvitz    Ece Kamar    Peter Lee    Yin Tat Lee    Yuanzhi Li    Scott Lundberg

Harsha Nori    Hamid Palangi    Marco Tulio Ribeiro    Yi Zhang

Microsoft Research

# The Hype

# Be realistic

# We are fully booked

# We are fully booked

# Candidate Selection

- Being here does **NOT** guarantee your acceptance

- We need to further select 20 students (10 from CS and 10 from LST)

- For CS student, revise your answer to our survey form

- For LST student, remember to fill out the from by **Apr 21, Friday 23:00**

- We will finish the selection process by **Apr 23, Sunday**

# Course overview

- Tue 12:15-13:45, 10 meetings (tentative, depends on our progress)

- Xudong Hong
    - xhong@coli.uni-saarland.de, C7 2, 2.02 by appointment.

- Ruitao Feng
    - fruitao@coli.uni-saarland.de, by appointment.

# Course overview

- Assessment:
  - Presentation of a chosen topic (paper)
  - Prepare questions
  - Participation in discussion
  - Peer review of weekly presentations
  - Report (for the 7 CP version)

# Prerequisites

- *How much do I need to know about neural networks/ NLP / CV to be able to follow the course?*

- 1. Students with practical experience with **NLP or CV** and want to
- deepen their theoretical understanding

- Or

- 2. Students who are familiar with neural NLP or CV (classes like NNIA,
- related seminars) and would like to gain practical experience

The most important prerequisite is you are really interested in the topic and you want to participate

# How familiar is "familiar"?

- Do you know what the following mean:
  - Seq2Seq
  - Self-attention
  - Transformer

If you're unclear about some, that's okay!
I will talk about them in the next meeting.

# Plan for today

1. **Introduction**
2. Reading List
3. Organization issues
4. Get to know each other

# What is Language and Vision (L&V)?

Language-vision learning is an interdisciplinary research area that involves the development of sophisticated algorithms and models that can comprehend and generate complex, context-rich information from both visual and linguistic inputs.

# **Why L&V Research?**

- This is how we interact with and learn from the world
  - Vision is a large portion of how humans perceive
  - Language is a large portion of how humans communicate
  - A smart Ai system should be able to perform well on both
  - One step towards AGI

# **Why L&V Research?**

- Multimodal learning has applications in various fields, including visual storytelling and human-computer interaction

Popular tasks:



Visual Question Answering
What color is the child's outfit?  Orange

Referring Expressions
child   sheep   basket   people sitting on chair

Multi-modal Verification
The child is petting a dog.  false

Caption-based Image Retrieval and *image captioning*
A child in orange clothes plays with sheep.

# Why L&V Research?

It is just fun!

- OpenAI DALL-E 2
  - creates realistic images from a given text prompt

# Common VL Tasks- Image Captioning

- Language generation
- Difficult automatic evaluation (BLEU, CIDEr, Rogue)

# Common VL Tasks- Visual Question Answering

- Elicit specific information from images
- Relatively easier evaluation (accuracy using string matching)



*Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering*, 2016

# Common VL Tasks- Visual Storytelling

- Creative Language generation



**Visually Grounded Story Generation**

*Jack was on a call with a client, getting stressed over a business deal that wasn't going well.*
*Jack put the phone down after an unsuccessful deal and decided to go get a coffee at the nearby coffee.*
*At the coffee shop, he started talking to the waiter Will about the unfortunate call.*
*Will told him he would convince the client to accept the deal if he could work for Jack.*
*Will then called the client and successfully struck the deal….*

# Seq2seq Models

# Seq2seq Models

Vinyals, O., Toshev, A., Bengio, S. and Erhan, D., 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 3156-3164).

# Plan for today

1. Introduction
2. **Reading List**
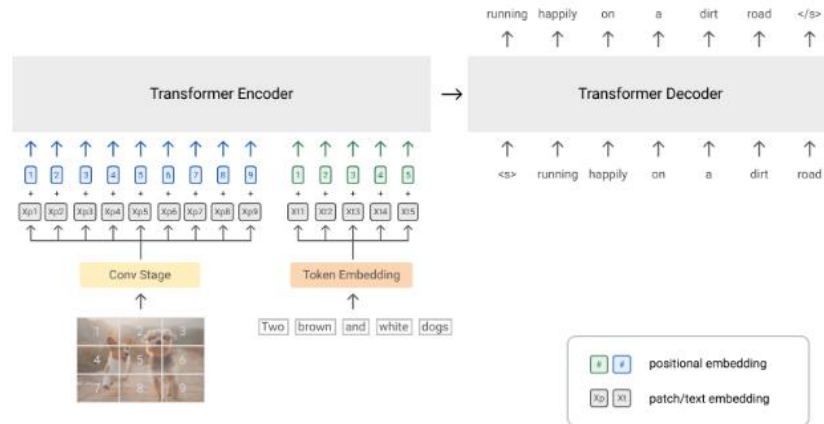3. Organization issues
4. Get to know each other

# Topics

- Pre-training – Image
    - Contrastive
    - Seq2Seq Pre-training – Image
- Pre-training – Video
- Multitask Learning
- Parameter Efficiency
    - Prompting
    - Prompt Tuning
    - Prefix-Tuning
    - Adapters
- Recent Generative Models
    - Text-to-Image
    - Diffusion
    - GPT
- Reinforcement Learning

# Seq2seq Pre-training - Image

- Wang, Z., Yu, J., Yu, A.W., Dai, Z., Tsvetkov, Y. and Cao, Y., SimVLM: Simple Visual Language Model Pretraining with Weak Supervision. In *International Conference on Learning Representations*.
- Singh, A., Hu, R., Goswami, V., Couairon, G., Galuba, W., Rohrbach, M. and Kiela, D., 2022. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 15638-15650).
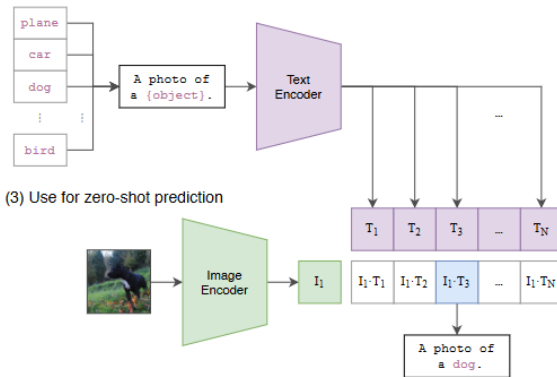
# Contrastive Pre-training - Image

- Li, J., Li, D., Xiong, C. and Hoi, S., 2022, June. [Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation](). In *International Conference on Machine Learning* (pp. 12888-12900). PMLR.
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J. and Krueger, G., 2021, July. [Learning transferable visual models from natural language supervision](). In International conference on machine learning (pp. 8748-8763). PMLR.
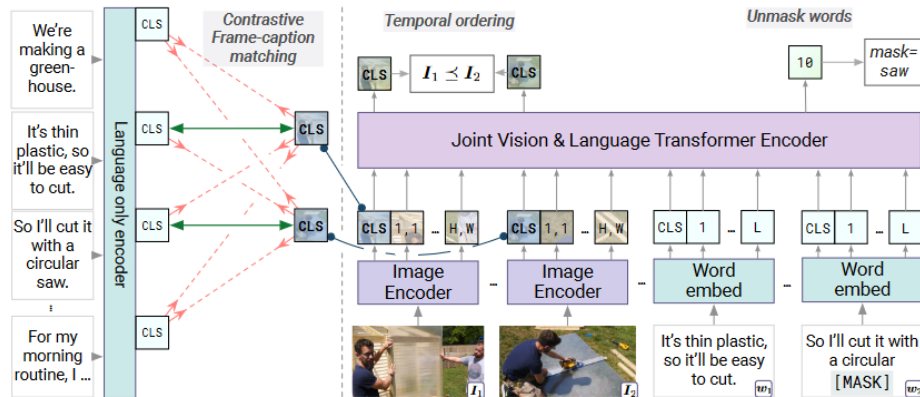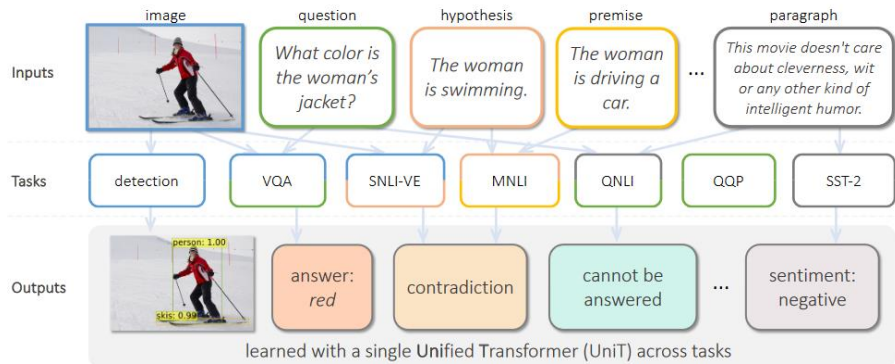
# Pre-training - Video

- Zellers, R., Lu, J., Lu, X., Yu, Y., Zhao, Y., Salehi, M., Kusupati, A., Hessel, J., Farhadi, A. and Choi, Y., 2022. Merlot reserve: Neural script knowledge through vision and language and sound. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 16375-16387).

- Zellers, R., Lu, X., Hessel, J., Yu, Y., Park, J.S., Cao, J., Farhadi, A. and Choi, Y., 2021. Merlot: Multimodal neural script knowledge models. *Advances in Neural Information Processing Systems*, *34*, pp.23634-23651.
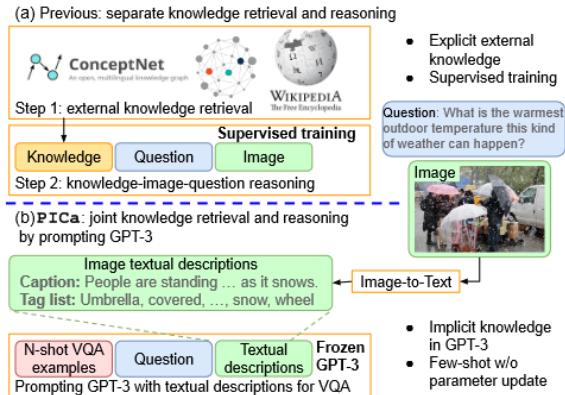
# Multitask Learning

- Cho, J., Lei, J., Tan, H. & Bansal, M.. (2021). Unifying Vision-and-Language Tasks via Text Generation. *Proceedings of the 38th International Conference on Machine Learning*, in *Proceedings of Machine Learning Research* 139:1931-1942

- Hu, R. and Singh, A., 2021. Unit: Multimodal multitask learning with a unified transformer. In Proceedings of the IEEE/CVF International Conference on Computer Vision (pp. 1439-1449).

- Wang, P., Yang, A., Men, R., Lin, J., Bai, S., Li, Z., Ma, J., Zhou, C., Zhou, J. and Yang, H., 2022, June. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning* (pp. 23318-23340). PMLR.
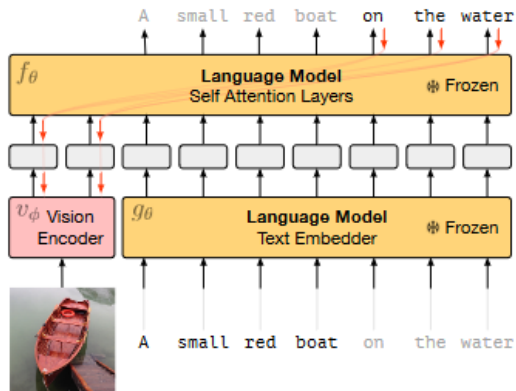
# Parameter Efficiency - Prompting

- Yang, Z., Gan, Z., Wang, J., Hu, X., Lu, Y., Liu, Z. and Wang, L., 2022, June. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 36, No. 3, pp. 3081-3089).
- Zeng, A., Wong, A., Welker, S., Choromanski, K., Tombari, F., Purohit, A., Ryoo, M., Sindhwani, V., Lee, J., Vanhoucke, V. and Florence, P., 2022. Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language. ICLR 2023. *arXiv e-prints*, pp.arXiv-2204.

# Parameter Efficiency - Prompt Tuning

- Tsimpoukelli, M., Menick, J.L., Cabi, S., Eslami, S.M., Vinyals, O. and Hill, F., 2021. Multimodal few-shot learning with frozen language models. Advances in Neural Information Processing Systems, 34, pp.200-212.
- Yu, Y., Chung, J., Yun, H., Kim, J. and Kim, G., 2021. Transitional adaptation of pretrained models for visual storytelling. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 12658-12668).
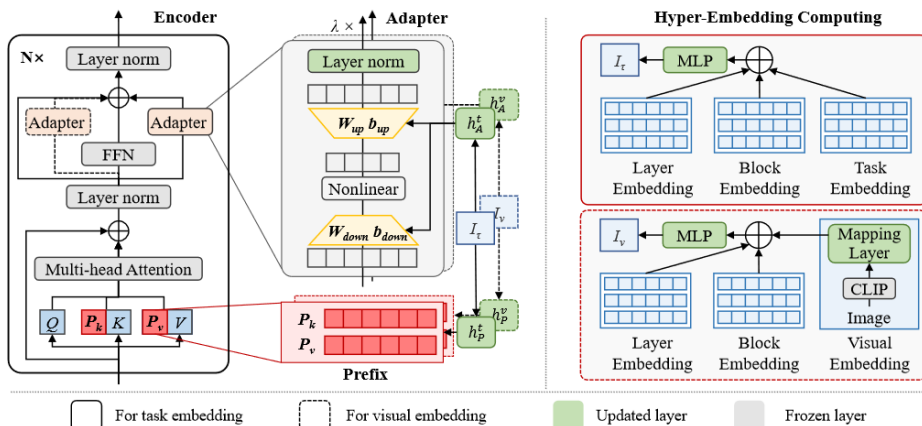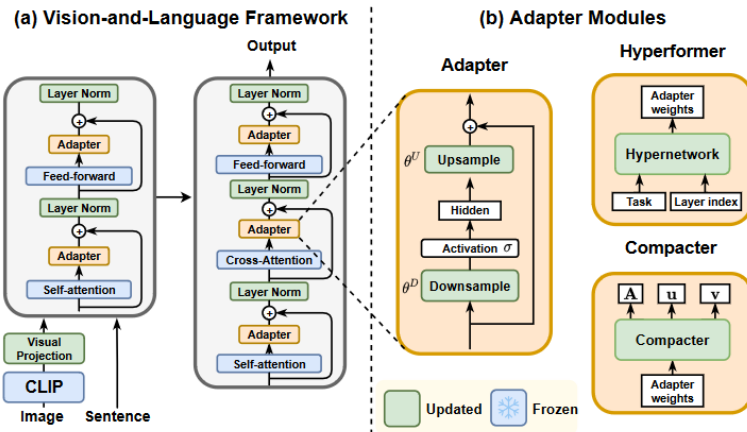
# Parameter Efficiency - Prefix-Tuning

- Zhang, Z., Guo, W., Meng, X., Wang, Y., Wang, Y., Jiang, X., Liu, Q. and Yang, Z., 2022. Hyperpelt: Unified parameter-efficient language model tuning for both language and vision-and-language tasks. arXiv preprint arXiv:2203.03878.
- Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B. and Lim, S.N., 2022, October. Visual Prompt Tuning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII* (pp. 709-727).
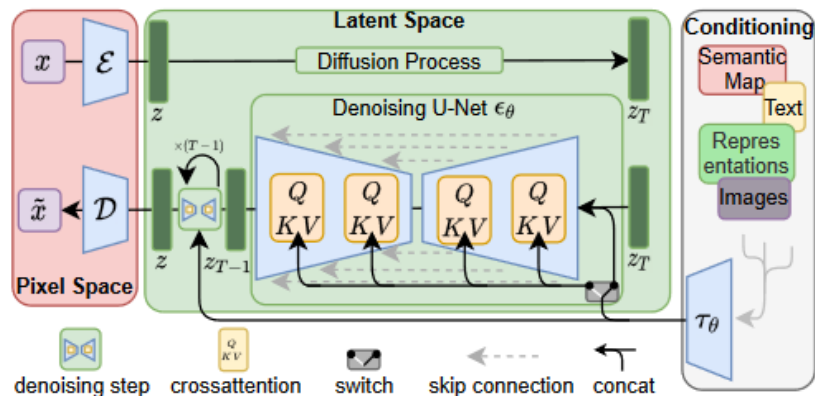
# Parameter Efficiency - Adapters

- Sung, Y.L., Cho, J. and Bansal, M., 2022. [Vl-adapter: Parameter-efficient transfer learning for vision-and-language tasks](#). In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 5227-5237).
- Sung, Y.L., Cho, J. and Bansal, M., 2022. [LST: Ladder Side-Tuning for Parameter and Memory Efficient Transfer Learning](#). In *Advances in Neural Information Processing Systems 2022*.

# Generative Model - Text-to-Image

- Rombach, R., Blattmann, A., Lorenz, D., Esser, P. and Ommer, B., 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 10684-10695).

# Generative Model - GPT

- OpenAI, 2023. GPT-4. Available at: https://openai.com/research/gpt-4. March 14, 2023.
  - (Optional) OpenAI (2023). GPT-4 Technical Report. ArXiv, abs/2303.08774.
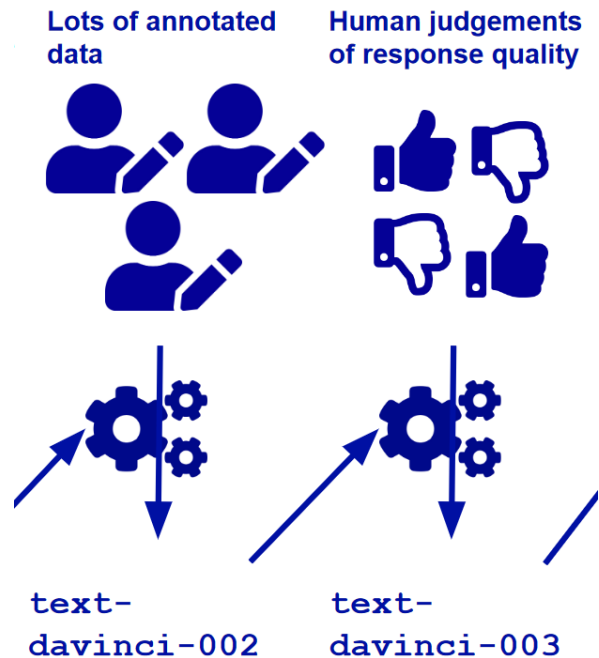
# Generative Model - GPT

- Wu, C., Yin, S., Qi, W., Wang, X., Tang, Z. and Duan, N., 2023. [Visual chatgpt: Talking, drawing and editing with visual foundation models](#). *arXiv preprint arXiv:2303.04671.*
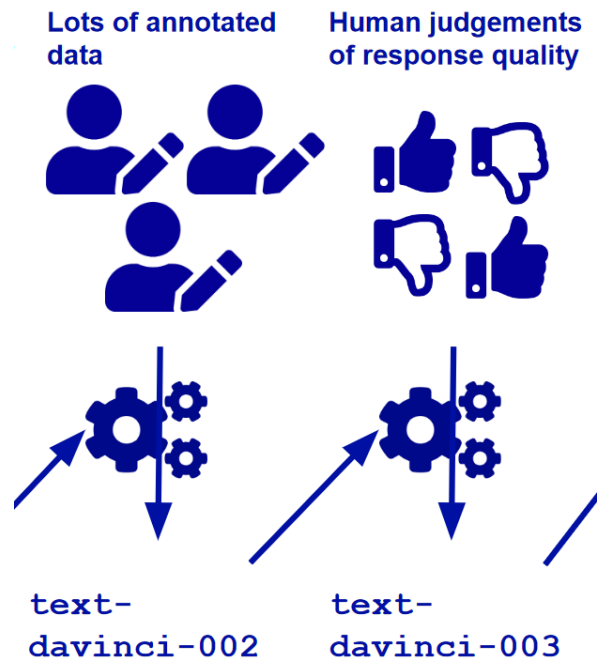
# Reinforcement Learning

- Why ChatGPT works?

# Reinforcement Learning

- Wang, X., Chen, W., Wang, Y.F. and Wang, W.Y., 2018, July. No Metrics Are Perfect: Adversarial Reward Learning for Visual Storytelling. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 899-909).

- Hu, J., Cheng, Y., Gan, Z., Liu, J., Gao, J. and Neubig, G., 2020, April. What makes a good story? designing composite rewards for visual storytelling. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 34, No. 05, pp. 7969-7976).

Lots of annotated data

Human judgements of response quality

text-davinci-002

text-davinci-003

# **Plan for today**

1. Introduction
2. Reading List
3. **Organization issues**
4. Get to know each other

# Peer review

- Helpful feedback for the presenters
- Single-blind review
- Don't worry. Your grades will not be affected by others' review on you. (On the contrary, the reviews will be graded.)

# **Quick guide for a good presentation**

Have a slide that answer these questions:

- Why is this work important?
  i.e. What is the general problem ?
  or what is the main question ?

- What is the purpose of this work?
  i.e. How does this work contribute to the sub-problem of the general problem?
  or what specific questions are this work trying to answer?

# Quick guide for a good presentation

- Problem 1, Problem 2, Problem 3
  Experiment 1, Experiment 2, Experiment 3,
  Results 1, Results 2, Results 3

  – Problem 1, Experiment 1, Results 1,
    Problem 2, Experiment 2, Results 2,
    Results 1, Results 2, Results 3

- "But I want to compare the experiments"

  – Just show Results 1 again alongside Results 2

Bad !!

Better !!

# **Quick guide for a good presentation**

Have a slide that answer these questions:

- Why is this work important?
  i.e. What is the general problem ?
  or what is the main question ?

- What is the purpose of this work?
  i.e. How does this work contribute to the sub-problem of the general problem?
  or what specific questions are this work trying to answer?

# Repeat:
# Practical arrangement

- For presenter:

    - Make an appointment with me and show me your slides by the **Wednesday** before your talk.

    - Build a demo with Google Colab

- For everyone else except the presenter:

    - Read the paper to be presented, and post questions on MS Teams by **Friday** before the talk.

    - Submit peer review of the presentation by **Friday** after the talk.

# Plan for today

1. Introduction
2. Reading List
3. Organization issues
4. **Get to know each other**

# It's your time!

Time to get to know each other

1. Your name and what you prefer to be called
2. Year in BSc/MSc
3. Native language(s)
4. Prior degree and areas of interest/specialization
5. Have you taken a related course before?



shutterstock.com · 1078542572

# References

- Yung, 2019, Discourse Relations: cognition, resources, NLP
- Iza, 2021, Controllable Text Generation

https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.videogamer.com%2Ftech%2Fai%2Fcould-gpt-5-have-artificial-general-intelligence%2F&psig=AOvVaw3v3Daodob39ixyZ-WHuWS8&ust=1681292172396000&source=images&cd=vfe&ved=2ahUKEwjpze7Gw6H-AhXjnCcCHXDjAQAQr4kDegUIARDNAQ

https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.youtube.com%2Fwatch%3Fv%3D1gThRMB7FRI&psig=AOvVaw3Vnh3StTvyod6LYlFe-PVe&ust=1681292130714000&source=images&cd=vfe&ved=2ahUKEwivw_6yw6H-AhW6XvEDHa5zAx4Qjhx6BAgAEA0

https://openai.com/research/gpt-4

xudonghong.me/mllv-uds