

Spatial Gesture Semantics

4. AI and Gesture Detection

Andy Lücking Alexander Henlein

Goethe University Frankfurt

July 28–August 01, 2025

Recap

Yesterday's lecture

- World-to-word direction of fit
- Classifier-based (computational) semantics
- Exemplification (extended exemplification)
- Informational evaluation heuristic

Today's lecture

- ML Primer: learning paradigms
- Building models like ChatGPT
- Multimodal foundations
- Gesture-detection pipeline
- Hands-on live demo
- Outlook & open questions

Introduction: Machine Learning, AI and Multimodality

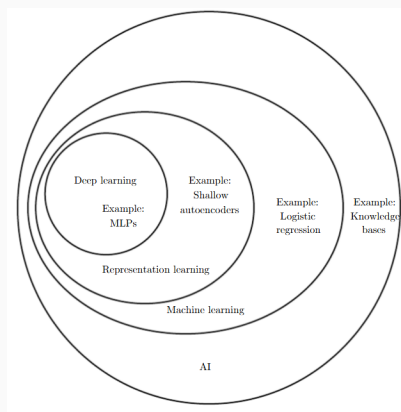
What is Machine Learning / AI?¹

Artificial Intelligence (AI)

- Umbrella term for techniques that enable machines to perform tasks we regard as “intelligent” (reasoning, perception, planning, language).

Machine Learning (ML)

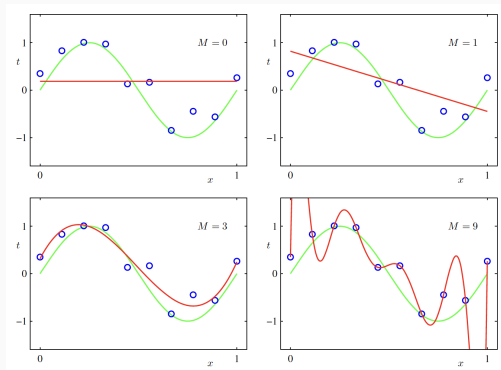
- Sub-field of AI: systems **learn** patterns from data instead of relying on hand-crafted rules.
- Core ingredients: large data \rightarrow model \rightarrow loss \rightarrow optimisation *rightarrow* evaluation.



¹ I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio (2016). **Deep learning**. Vol. 1. MIT press Cambridge

What is Machine Learning (ML)?³

- Instead of writing explicit rules, ML finds **patterns** in data.
- At its core: ML = fitting a function to data.
- Useful when the rules are too complex, fuzzy, or unknown - e.g., how gestures vary across speakers and contexts.

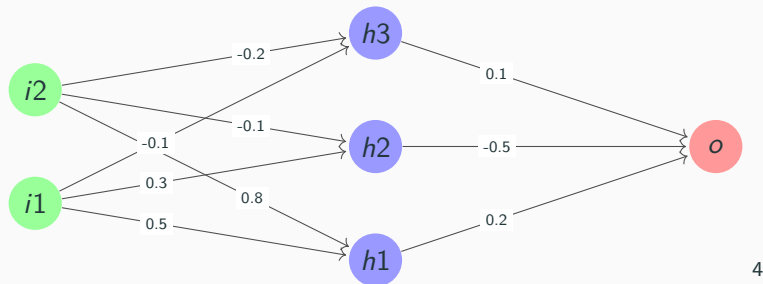


²https://amitrajan012.github.io/post/pattern-recognition-chapter-1-introduction_1/

³ C. M. Bishop and N. M. Nasrabadi (2006). **Pattern recognition and machine learning**. Vol. 4. Springer

How are Neural Networks trained?

Is the weather suitable for picnics?

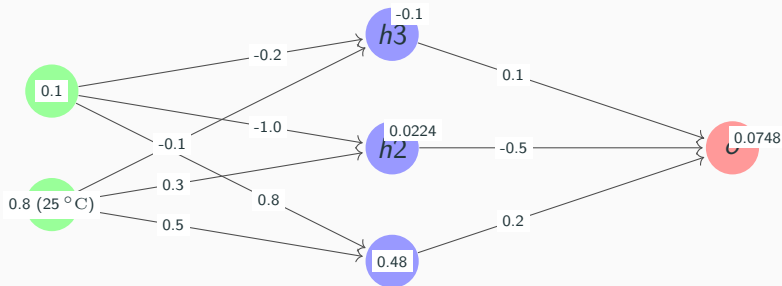


$i1$: temperature
 $i2$: risk of rain
 o : picnic score

⁴Template: <https://tikz.net/regular-vs-bayes-nn/>

How are Neural Networks trained? - Example Input

Is the weather suitable for picnics?



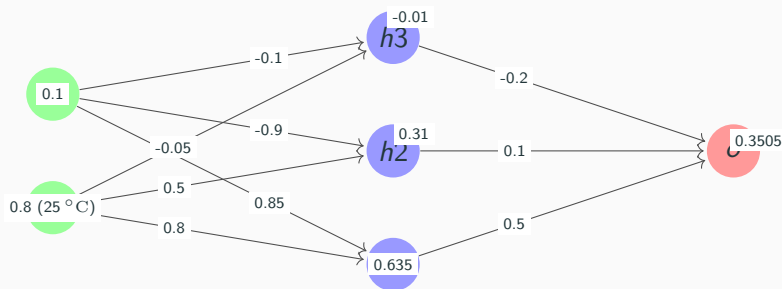
i1: temperature
i2: risk of rain
o: picnic score

5

⁵Template: <https://tikz.net/regular-vs-bayes-nn/>

How are Neural Networks trained? - Backpropagation

Is the weather suitable for picnics?



i1: temperature
i2: risk of rain
o: picnic score

6

⁶Template: <https://tikz.net/regular-vs-bayes-nn/>

Learning Paradigms in ML⁷

- Supervised → Learn to predict labels.
- Unsupervised → Find structure or clusters.
- Self-Supervised → Predict part of data from other parts.
- Semi-Supervised → Leverage a few labels with lots of unlabeled data.
- Reinforcement → Learn good decisions over time.

⁷ C. M. Bishop and N. M. Nasrabadi (2006). [Pattern recognition and machine learning](#). Vol. 4. Springer; V. Rani et al. (2023). “Self-supervised learning: A succinct review”. In: [Archives of Computational Methods in Engineering](#) 30, 2761–2775

What ML Can (and Can't) Do for Us⁸

Can do

- **Detect** classes even from noisy data.
- **Cluster** and quantify variation.
- **Learn** useful representations from raw data.
- **Support** large-scale studies of form and use.

Can't do

- **Understand** meaning on its own.
- **Replace** semantic theory or manual insight.
- **Handle** open-ended or subtle communicative functions (yet).
- **Guarantee** fairness, explainability, or trustworthiness out of the box.

⁸ E. M. Bender and A. Koller (2020). “Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data”. In: [Proc. of the 58th Annual Meeting of the Association for Computational Linguistics](#), 5185–5198; G. Marcus and E. Davis (2019). [Rebooting AI: Building artificial intelligence we can trust](#). Vintage

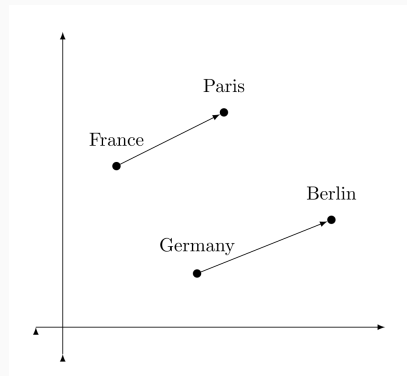
Embeddings: Representing Data as Vectors⁹

What are Embeddings?

- Continuous vector representation of discrete items (words, tokens, images).
- Geometric proximity \Leftrightarrow semantic similarity.

Why Important for LLMs

- Input tokens mapped to embeddings learned during training.
- Enable efficient dot-products, generalisation, and transfer across tasks.

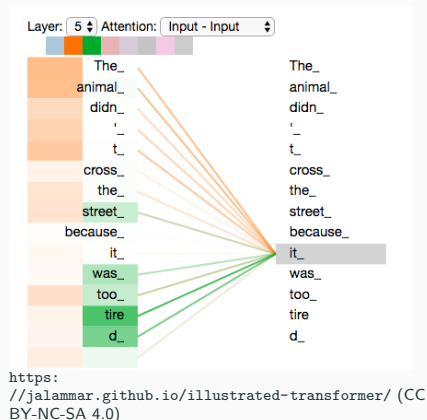


(Wikimedia Foundation, Inc. Original up- loader was Cbarr (WMF), CC BY-SA 3.0, File:RobGrindes-shrug-143px.png)

⁹ T. Mikolov et al. (2013). “Distributed Representations of Words and Phrases and their Compositionality”. In: [Advances in Neural Information Processing Systems](#); T. Mikolov, K. Chen, G. Corrado, and J. Dean (2013). “Efficient Estimation of Word Representations in Vector Space”. In: [1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings](#)

Attention Mechanism¹⁰

- Dynamically weights input elements based on relevance.
- **Self-attention**: queries, keys, values from same sequence.
- **Multi-head**: parallel views capture diverse relations.
- Powers the Transformer architecture and modern LLMs.



¹⁰ A. Vaswani et al. (2017). "Attention is All you Need". In: [Advances in Neural Information Processing Systems](#)

(Multimodal) Large Language Models

How to train my own ChatGPT¹²

1. Data Collection & Preprocessing → Clean, filter, deduplicate, normalize, tokenize.
2. (Self-supervised) Pretraining → Next-token prediction.
3. Post-Training → Reinforcement Learning from Human Feedback (RLHF)¹¹.
4. Evaluation → For performance, safety, bias, hallucination.
5. Deployment & Iteration → Frequent monitoring and updated.

¹¹ L. Ouyang et al. (2022). “Training language models to follow instructions with human feedback”. In: [Proc. of the 36th International Conference on Neural Information Processing Systems](#)

¹² OpenAI et al. (2024). [GPT-4 Technical Report](#). arXiv: 2303.08774 [cs.CL]

Step 1: Data Collection & Preprocessing¹³

Goal:

Prepare high-quality, diverse input for training.

Sources:

- Web text
- Books, Wikipedia
- Forums, code repositories
- Internal/proprietary data

- **Filtering**: remove low-quality, toxic, or irrelevant content.
- **Deduplication**: avoid overfitting to repeated content.
- **Normalization**: standardize text (e.g., lowercase, punctuation).
- **Tokenization**: convert text into input tokens.
- **Balancing**: ensure coverage across domains (e.g., code vs. dialogue).

¹³ L. Gao et al. (2020). [The Pile: An 800GB Dataset of Diverse Text for Language Modeling](#). arXiv: 2101.00027 [cs.CL]

Excursus: What is Tokenization in the Context of LLMs?¹⁴

Goal:

Convert raw text into units the model can understand.

Why not characters or words?

- Characters: too granular, inefficient
- Words: ambiguous, too many
- Tokens: trade-off

- Use subword units (e.g. “play“, “#ing“; “un“, “#believable“).
- Based on algorithms like **Byte-Pair Encoding (BPE)** or **Unigram LM**.
- Allows handling of rare and unknown words.
- Example: "I really enjoyed my time in Bochum." → ["I", "really", "enjoy", "#ed", "my", "time", "in", "Boch", "#um", "."]

¹⁴https://huggingface.co/docs/transformers/tokenizer_summary

Step 2: (Self-supervised) Pretraining

Goal:

Teach the model general language understanding.

Method:

- Predict next token
- No human labels needed
- Very large dataset

- Objective: $P(\text{token}_t \mid \text{token}_{1..t-1})$
- Transformer architecture (e.g. decoder-only).
- Trained on trillions of tokens.
- Requires massive compute (TPUs, GPUs).
- Learns grammar, facts, reasoning, coding patterns.

Step 3: Post-Training (Alignment)

Goal:

Make the model helpful, safe, and aligned with human values.

Steps:

- Supervised fine-tuning (SFT)
- RLHF (Reinforcement Learning from Human Feedback)¹⁵

- Human-written prompt-response pairs.
- Rank model outputs → train a reward model.
- Fine-tune the base model using Reinforcement Learning.
- Encourages helpful and non-toxic responses.
- Aligns model with human intent.

¹⁵ Y. Bai et al. (2022). [Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback](#). arXiv: 2204.05862 [cs.CL]

Step 4: Evaluation¹⁶

Goal:

Assess model quality, safety, and behavior before release.

Types:

- Quantitative tests
 - Human evaluations
 - Red-teaming
- Benchmarking (MMLU, HellaSwag, etc.).
 - Prompt diversity testing and edge cases.
 - Detect bias, toxicity, hallucinations.
 - Internal and external safety audits.
 - Analyze model confidence and calibration.

¹⁶ Y. Chang et al. (2024). “A survey on evaluation of large language models”. In: [ACM transactions on intelligent systems and technology](#) 15, 1–45

Step 5: Deployment & Iteration

Goal:

Safely deploy the model and keep improving it through usage.

Cycle:

- Launch → Monitor → Improve
- Continuous feedback loop

- Model exposed via APIs, apps (e.g. ChatGPT).
- Usage analytics + human feedback collected.
- Updates: bugfixes, safety patches, new features.
- Ongoing fine-tuning and A/B testing.
- Data pipeline refinement based on usage.

Step 6a: Image Encoder¹⁷

Goal:

Convert an image into a vector representation (embeddings).

Common Encoders:

- CLIP (ViT)
- ResNet
- SigLIP
- Vision Transformer (ViT)

- **Input:** raw image pixels
- **Output:** sequence of image embeddings (like tokens)
- Pretrained on image-text pairs (e.g., from web)
- Encoded images are fed into the language model as part of the prompt
- Can capture visual objects, layout, and spatial info

¹⁷ A. Radford et al. (2021). “Learning Transferable Visual Models From Natural Language Supervision”. In: [CoRR](#) abs/2103.00020. arXiv: 2103.00020

Step 6b: Aligning Modalities¹⁸

Goal:

Bridge the gap between visual and textual representations.

Why align?

- Image + text are from different distributions
- Need unified input for Transformer

- **Projection Layer:** maps image embeddings to LLM token space
- **Concatenation:** image embeddings placed before or between text tokens
- **Joint Training:** learn to ground vision in language tasks
- Enables multimodal reasoning, captioning, and VQA

¹⁸ H. Liu, C. Li, Q. Wu, and Y. J. Lee (2023). “Visual instruction tuning”. In: [Advances in neural information processing systems](#) 36, 34892–34916

Step 6c: Multimodal Training Tasks¹⁹

Goal:

Teach the model to understand and reason over image-text pairs.

Common Task Types:

- Image captioning
- Visual question answering (VQA)
- OCR + scene text recognition
- Referring expression resolution

- **Image → Text**: Generate captions or summaries
- **Image + Text → Text**: Answer questions about the image
- Use instruction-following prompts: "Describe this image.", "Where is the cat?"
- Supervised training followed by instruction tuning
- Important for grounding language in perception

¹⁹ Y. Zhang et al. (2024). [LLaVAR: Enhanced Visual Instruction Tuning for Text-Rich Image Understanding](#). arXiv: 2306.17107 [cs.CV]

Excursus: Let's take a look at this type of data set.

`http://captions.christoph-schuhmann.de/eval_laion/eval.html`
`https://laion.ai/projects/`

Scaling and Advanced Concepts

Goal:

Push efficiency, capability, and scalability of large models.

Key Concepts:

- Mixture of Experts (MoE)
- Sparse Attention
- Retrieval-Augmented Generation (RAG)
- Instruction Tuning

- **MoE**: Only a subset of model components (experts) is activated per input. Greatly reduces compute cost while increasing parameter count.
- **Sparse Attention**: Improves efficiency in very long-context models. Models learn to focus selectively.
- **RAG**: Combines LLMs with external search/indexes. Augments generation with real-world knowledge.
- **Instruction Tuning**: Further trains models to follow natural language commands more reliably. Key to usability.

Why shrink models?

- Lower inference latency and energy cost.
- Fit on-device / edge hardware.
- Enable private, offline use.
- Reduce carbon footprint.

Main techniques

- **Quantization**: 8-bit/4-bit weights.
- **Pruning**: remove redundant weights or neurons.
- **Distillation**: train a smaller student on teacher outputs.
- **PEFT**: LoRA, Adapters, ...

²⁰ S. Park, J. Choi, S. Lee, and U. Kang (2024). [A Comprehensive Survey of Compression Algorithms for Language Models](#). arXiv: 2401.15347 [cs.CL]

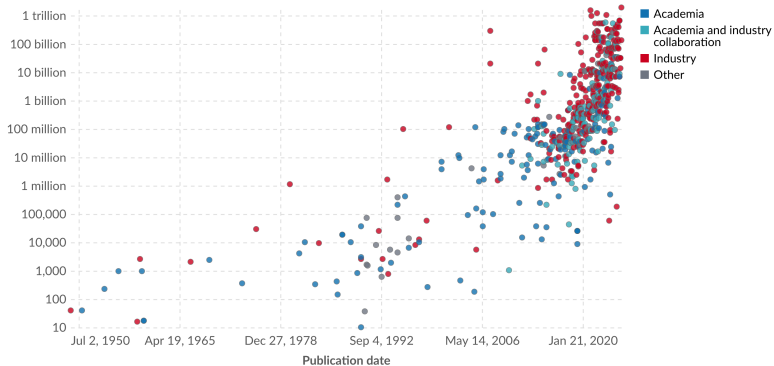
Where is this going?²¹

Parameters in notable artificial intelligence systems

Our World
in Data

Parameters are variables in an AI system whose values are adjusted during training to establish how input data gets transformed into the desired output; for example, the connection weights in an artificial neural network.

Number of parameters



Data source: Epoch (2025)

OurWorldinData.org/artificial-intelligence | CC BY

Note: Parameters are estimated based on published results in the AI literature and come with some uncertainty. The authors expect the estimates to be correct within a factor of 10.

²¹<https://ourworldindata.org/grapher/artificial-intelligence-parameter-count>

Where is this going?²²

Major Large Language Models (LLMs)

ranked by capabilities, sized by billion parameters used for training

CLICK LEGEND ITEMS TO FILTER

anthropic chinese google meta microsoft mistral openAI other

Parameters (Bn) open access

search...

show only: all

100 MMLU

89.8 = human expert

80

70+ IDEAL

60

40

20

pre-2022

2022

2023

2024

2025

David McCandless, Tom Evans, Paul Barton
informationisbeautiful // Jan 2024

MMLU = benchmark for measuring LLM capabilities
* = parameters undisclosed // source: LifeArchitecture // data

²²<https://informationisbeautiful.net/visualizations/the-rise-of-generative-ai-large-language-models-llms-like-chatgpt/>

Critical Reflections on Large Language Models

Scaling is powerful, but not free

Open Issues:

- Environmental cost
- Evaluation transparency
- Data ethics
- Model accessibility

- **Compute cost:** training GPT-3 used hundreds of PFLOPs-days; carbon footprint estimated at hundreds of tons CO²³.
- **Financial cost:** GPT-4-level training estimated at millions of dollars; access to compute increasingly centralized²⁴.
- **Data concerns:** Training data scraped from the web—raises copyright, consent, and fairness issues²⁵.
- **Evaluation gaps:** Benchmarks often narrow and do not capture robustness, fairness, or real-world alignment²⁶.

²³<https://news.mit.edu/2025/explained-generative-ai-environmental-impact-0117>

²⁴<https://hai.stanford.edu/ai-index/2025-ai-index-report> (page 65)

²⁵ P. Samuelson (2023). “Generative AI meets copyright”. In: *Science* 381, 158–161

²⁶ T. R. McIntosh et al. (2025). “Inadequacies of large language model benchmarks in the era of

Gesture Detection

Step 1: Video Input

Input Type:

- Usually RGB frames or video stream
- Optionally with depth or IR

Preprocessing:

- Resize, normalize
- Frame extraction or windowing
- Optional face/hands segmentation

Challenges:

- Varying lighting and backgrounds
- Occlusion (e.g. hands crossing)
- Real-time constraints (latency, FPS)
- Device variability (camera quality)

Step 2: Hand Pose Detection

Goal: Localize key hand joints
(2D/3D)

- Wrist, knuckles, fingertips
- Input: video frame or cropped hand region
- Accuracy depends on input quality and occlusions
- Trade-off between model size and speed
- 3D pose enables gesture generalization (rotation invariance)
- Tracking is often fused with detection for consistency

Top-Down (two-step)

- Detect each **person/hand** first (e.g. bounding box).
- Run a pose/gesture network **inside** every box.
- **Pros:** high single-instance accuracy; leverages powerful object detectors.
- **Cons:** time scales $\#$ people; errors cascade from detector;

Bottom-Up (part-based)

- Detect **all keypoints** in the frame at once.
- Group points into individuals via part-affinity / clustering.
- **Pros:** cost nearly constant to crowd size; robust to missed boxes.
- **Cons:** grouping step can fail in heavy occlusion; slightly lower peak accuracy.

²⁷ R. Yue, Z. Tian, and S. Du (2022). “Action recognition based on RGB and skeleton data sets: A survey”. In: **Neurocomputing** 512, 287–306

Excursus: Common Models

- **MediaPipe**²⁸ (**Google**) - Lightweight, real-time framework for hand pose (21 keypoints per hand). Ideal for single-person tracking on mobile and web. Integrated into many apps and easy to use.
- **OpenPose**²⁹ (**CMU**) - Pose model supporting hands, body, and face. Requires GPU. Strong multi-person support. Still a popular baseline in research.
- **MMPose**³⁰ (**OpenMMLab**) - Modular PyTorch framework supporting many backbones and datasets. Includes whole-body hand keypoints and supports both 2D and 3D models. Great for custom experiments.
- **Sapiens**³¹ (**Meta**) - Newest, high-resolution foundation model with 308 keypoints (including hands). Designed for detailed, frame-by-frame offline analysis, not real-time use.

²⁸<https://github.com/google-ai-edge/mediapipe>

²⁹<https://github.com/CMU-Perceptual-Computing-Lab/openpose>

³⁰<https://github.com/open-mmlab/mmpose>

³¹<https://github.com/facebookresearch/sapiens>

Step 3: Feature Encoding³²

Goal: Convert hand pose data into useful input features for ML models.

- Raw 2D/3D keypoints
- Distances between joints
- Angles between fingers
- Motion vectors (velocity, acceleration)

Feature Types:

- Frame-based (pose snapshot)
- Sequence-based (temporal movement)
- Hand-crafted vs. deep-learned embeddings
- Normalize for translation, scale, rotation

³² P. Molchanov et al. (2016). "Online Detection and Classification of Dynamic Hand Gestures with Recurrent 3D Convolutional Neural Networks". In: [2016 IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#), 4207–4215

Step 4: ML Model Training³³

Input: Encoded features or pose sequences

Common Models:

- Classical ML: SVM, k-NN, Random Forest
- Deep Learning: MLP, CNN, LSTM, Transformers
- Spatio-temporal models for gesture dynamics

Training Considerations:

- Supervised learning with gesture labels.
- Augment data for generalization.
- Cross-subject and cross-session robustness.

³³ J. J. Ojeda-Castelo, M. d. L. M. Capobianco-Uriarte, J. A. Piedra-Fernandez, and R. Ayala (2022). “A survey on intelligent gesture recognition techniques”. In: [IEEE Access](#) 10, 87135–87156

Step 5: Gesture Prediction³⁴

Goal: Classify or detect user gestures in real-time.

Output Types:

- Static: e.g., "Thumbs up", "Open hand"
- Dynamic: e.g., "Swipe left", "Draw circle"

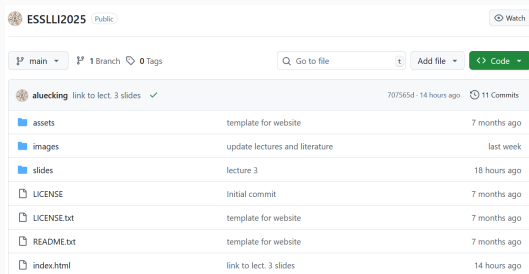
Deployment:

- Smooth output with tracking or temporal smoothing.
- Handle uncertain input with confidence thresholds.

³⁴ P. Molchanov et al. (2016). "Online Detection and Classification of Dynamic Hand Gestures with Recurrent 3D Convolutional Neural Networks". In: [2016 IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#), 4207–4215

Hands On Example

EnvisionBOX (https://github.com/aluecking/ESSLI2025)



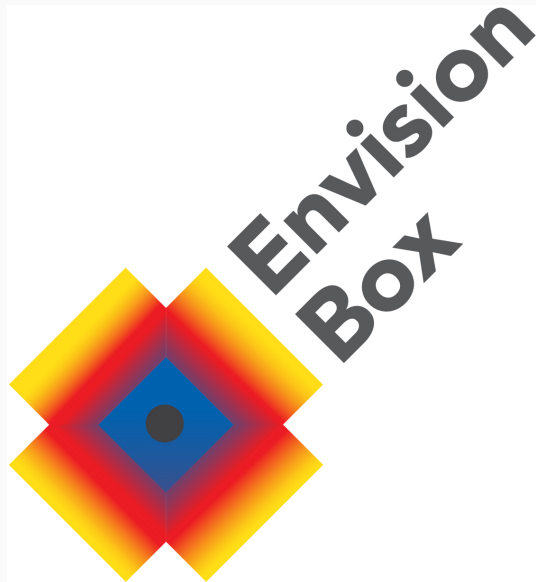
Public

Watch

main 1 Branch 0 Tags

Go to file Add file <> Code

aluecking link to lect. 3 slides ✓		707565d - 14 hours ago	11 Commits
assets	template for website	7 months ago	
images	update lectures and literature	last week	
slides	lecture 3	18 hours ago	
LICENSE	initial commit	7 months ago	
LICENSE.txt	template for website	7 months ago	
README.txt	template for website	7 months ago	
index.html	link to lect. 3 slides	14 hours ago	



Future Directions

Where are the LLMs for Gesture Detection?³⁶

Data & Representation

- **Sparse paired data:** very few corpora link *gesture key-points + language*; most instruction sets only contain captions.
- **Temporal mismatch:** LLMs digest static images; gestures are >30 FPS sequences \Rightarrow token explosion.
- **Modality gap:** 2-D RGB misses depth, skeleton cues essential for fine-grained hand motion.

³⁵ M. Cai et al. (2024). “Making Large Multimodal Models Understand Arbitrary Visual Prompts”. In: [IEEE Conference on Computer Vision and Pattern Recognition](#)

³⁶ D. Feng et al. (2025). “PoseLLaVA: Pose Centric Multimodal LLM for Fine-Grained 3D Pose Manipulation”. In: [Proceedings of the AAAI Conference on Artificial Intelligence](#) 39, 2951–2959; D. Zhang, T. Hussain, W. An, and H. Shouno (2025). [LLaVA-Pose: Enhancing Human Pose and Action Understanding via Keypoint-Integrated Instruction Tuning](#). arXiv: 2506.21317 [cs.CV]

Model Capabilities

- **Spatial precision:** overlay-based prompting (e.g. ViP-LLaVA) leaves 3-D joint reasoning unsolved.³⁵
- **Reasoning granularity:** current MLLMs excel at object semantics, but struggle with fine motor actions (pinch, swipe).
- **Safety/bias:** ambiguous gestures vary culturally; no robust alignment or policy-tuning yet.

Where are the LLMs for Gesture Detection?

PoseLLaVA³⁷

- **Model changes:** adds a pose encoder plus a cross-attention into LLaVA for global & local pose-image.
- **Created datasets:** adds **PosePart** (135 K single-body-part triplets) and combines Human3.6M (300 K), PoseScript (100 K) and PoseFix (135 K).
- **Finetuning task:** three-stage pipeline: pose-image contrastive pre-align, LLM pre-train on pose generation, unified instruction-tuning over estimation / generation / adjustment.

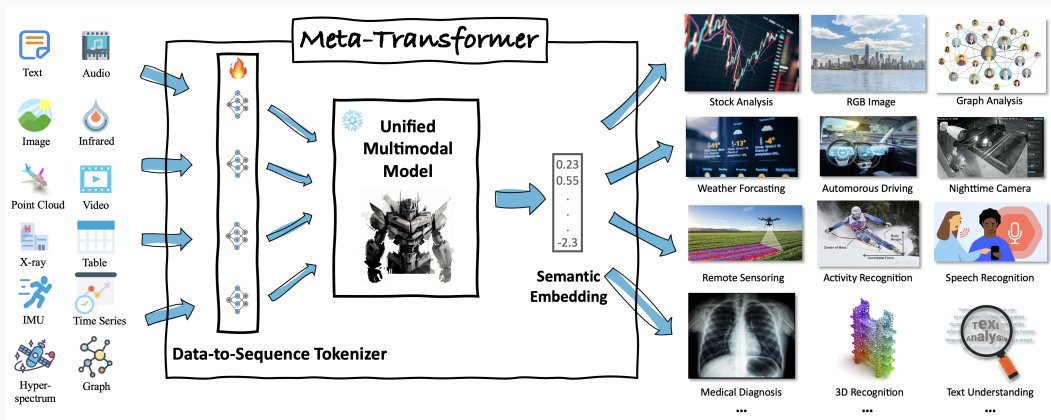
³⁷ D. Feng et al. (2025). "PoseLLaVA: Pose Centric Multimodal LLM for Fine-Grained 3D Pose Manipulation". In: [Proceedings of the AAAI Conference on Artificial Intelligence](#) 39, 2951–2959

LLaVA-Pose³⁸

- **Model changes:** no architectural edits - retains LLaVA-1.5 and simply augments prompts with 2-D keypoints.
- **Created datasets:** auto-generates 200328 COCO-based keypoint-aware instructions (conversation / description / reasoning); also publishes the **E-HPAUB** benchmark for evaluation.
- **Finetuning task:** full-model fine-tune for one epoch; objective is richer chat, description & reasoning about human pose/action scenes.

³⁸ D. Zhang, T. Hussain, W. An, and H. Shouno (2025). [LLaVA-Pose: Enhancing Human Pose and](#)

Concept: Meta-Transformer³⁹



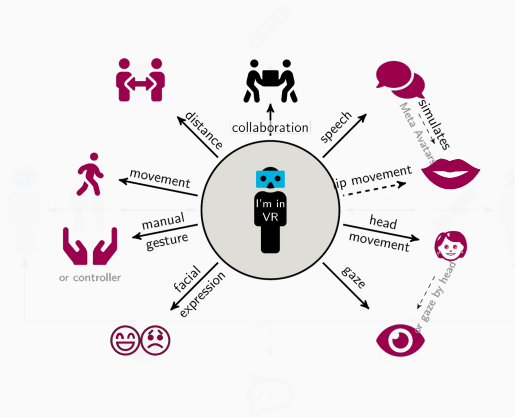
³⁹ Y. Zhang et al. (2023). “Meta-transformer: A unified framework for multimodal learning”. In: [arXiv preprint arXiv:2307.10802](https://arxiv.org/abs/2307.10802)

Why most LLMs stay text & image⁴⁰

- **Compute / memory** - universal tokens make long sequences; self-attention costly.
- **Limited generative** - good at unimodal perception, unclear for cross-modal generation.
- **Dataset gaps** - few richly multimodal pairs to unlock full promise.
- **Data scale** - trillions of web tokens, billions of captions; far fewer paired corpora.
- **Token explosion** - 10s of 30FPS video \approx 900 frames \Rightarrow hundreds of tokens.
- **ROI focus** - chat, code, doc QA, image help already monetize; niche sensors give uncertain payoff.
- **Tooling maturity** - CLIP/LLaVA pipelines are production-ready; multimodal 3-D/audio stacks still research-grade.

⁴⁰ A. Henlein et al. (2024). “An Outlook for AI Innovation in Multimodal Communication Research”. In: [Digital Human Modeling and Applications in Health, Safety, Ergonomics and Risk Management](#), 182–234; J. Jiang et al. (2025). [Token-Efficient Long Video Understanding for Multimodal LLMs](#). arXiv: 2503.04130 [cs.CV]; A. Kumar, M. M. Salim, D. Camacho, and J. H. Park (2025). “A comprehensive survey on large language models for multimedia data security: challenges and solutions”. In: [Computer Networks](#) 267, 111379; Y. Zhang et al. (2023). “Meta-transformer: A unified framework for multimodal learning”. In: [arXiv preprint arXiv:2307.10802](#)

- A VR-based simulation system.
- Multi-user collaborative tool.
- Users are represented by Meta Avatars.



⁴¹ A. Mehler et al. (2023). "A Multimodal Data Model for Simulation-Based Learning with Va.Si.Li-Lab". In: [Digital Human Modeling and Applications in Health, Safety, Ergonomics and Risk Management](#), 539–565

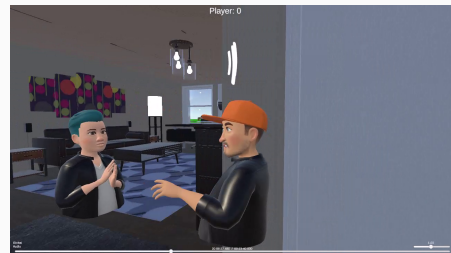




⁴² A. Lücking et al. (2010). “The Bielefeld Speech and Gesture Alignment Corpus (SaGA)”. In: *Multimodal Corpora: Advances in Capturing, Coding and Analyzing Multimodality*. 7th International Conference for Language Resources and Evaluation, 92–98

73 dialogues involving 146 speakers.

	Speaking time	# Tokens
total:	12:44:37	92,923
Router:	8:17:19	70,517
Follower:	4:27:18	22,406
Avg. Router:	0:06:49	1,273
Avg. Follower:	0:03:40	966
Avg. Dialogue:	0:10:28	307



⁴³ Lücking, Voll, Rott, Henlein, Mehler (2025). “Head and hand movements during turn transitions: data-based multimodal analysis using the Frankfurt VR Gesture–Speech Alignment Corpus (FraGA)”. In: accepted. 29th edition of the SemDial workshop series