# *From ChatGPT to GAIA-1: On Generative Sequence Models in Speech, Language, and Vision*

Tim Fingscheidt, Timo Lohrenz, Zhengyang Li, Björn Möller

# GAIA-1
## New possibilities in autonomous driving R&D

**What can GAIA-1 do?**

- It allows multimodal prompting with a video, text, and action, and …
- … hallucinates a realistic continuation of the video, under text and action constraints

**Why is GAIA-1 interesting** for autonomous driving?

- *Offline*: Generating „unlimited" video training/validation data, including some corner cases not seen in GAIA-1 training material
- *Online*: Can it even provide „a number of futures" for better trajectory planning?

**How does GAIA-1 technically work?**

- GAIA-1 is a generative sequence world model for autonomous driving R&D
- Any video, text, and action prompts are individually tokenized
- After tokenization, the prompts are conditions into a recurrently excecuted world model …
- … which delivers a future image token sequence, …
- … which is input to a recurrently executed diffusion video decoder, delivering a respective video sequence from it.

## GAIA-1:
## A Generative World Model for Autonomous Driving

**Anthony Hu**[*]  **Lloyd Russell**[*]  **Hudson Yeo**[*]  **Zak Murez**  **George Fedoseev**

**Alex Kendall**  **Jamie Shotton**

[1] Hu, Anthony, et al. "GAIA-1: A Generative World Model for Autonomous Driving." *arXiv preprint arXiv:2309.17080* (2023).

17 June 2023 | Research

WAYVE

## Introducing GAIA-1: A Cutting-Edge Generative AI Model for Autonomy

Technische Universität Braunschweig

ifN
Institut für Nachrichtentechnik

[Hu, Anthony, et al. "GAIA-1: A Generative World Model for Autonomous Driving." *arXiv preprint arXiv:2309.17080* (2023)]

https://youtu.be/5Jx2QgEUZUI
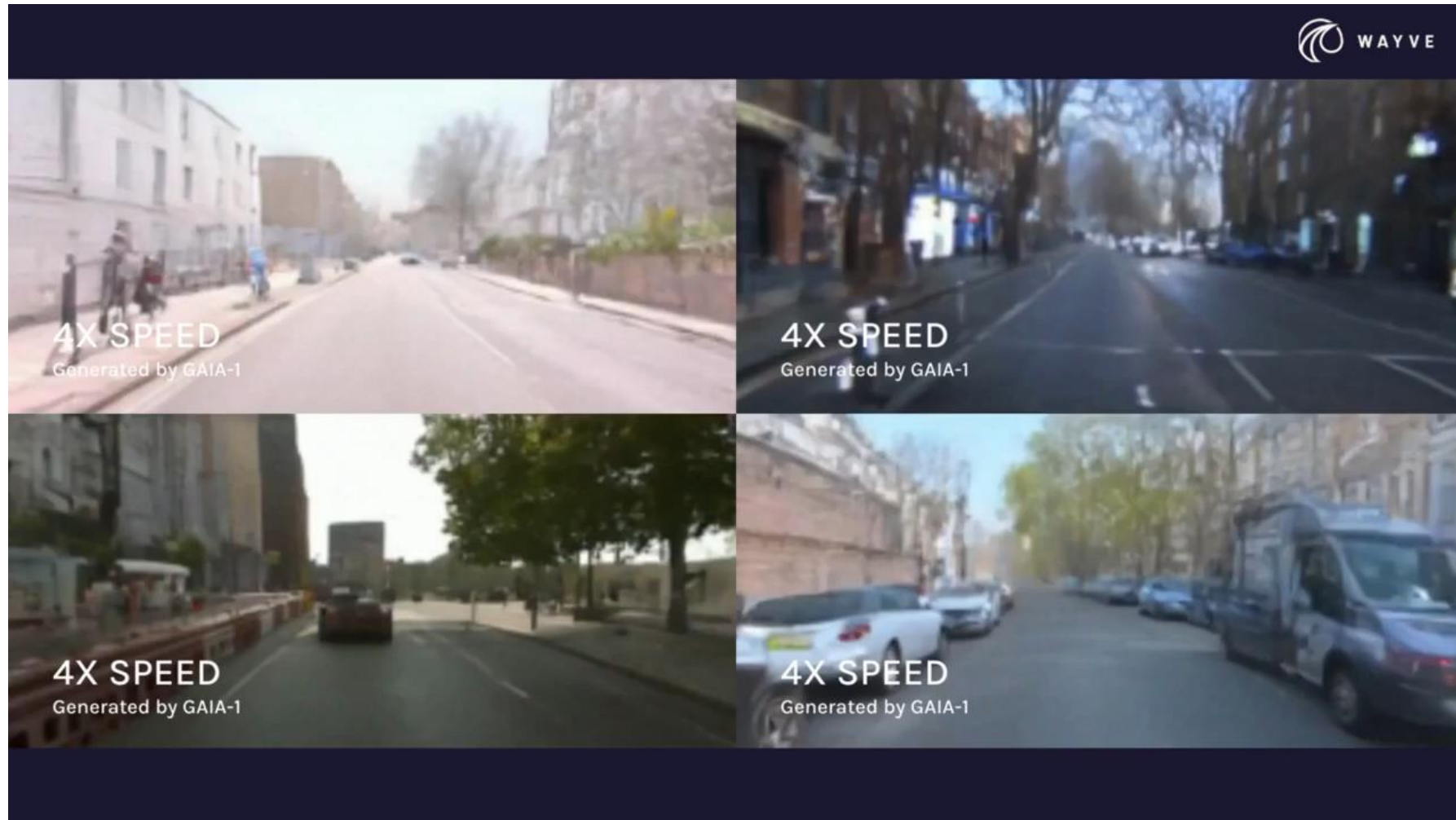
From ChatGPT to GAIA-1: On Generative Sequence Models in Speech, Language, and Vision  |  Prof. Tim Fingscheidt  |  DEC 6, 2023  |  4

GAIA-1 technology observations:

- Text tokenizer, world model, and video decoder are … transformer encoder and/or decoder models

  … executed recurrently to produce output step-by-step

- The very same recurrent execution of transformer models is also used in …

  … end-to-end automatic speech recognition (E2E ASR): ENC/DEC transformer, and in …

  … large language models (LLMs, e.g., ChatGPT): DEC transformer

> The formulation of the world modeling task in GAIA-1 shares a commonality with the approach frequently used in large language models (LLMs). In both instances, the task is streamlined to focus on predicting the next token. Although this approach is adapted for world modeling in GAIA-1 rather than the traditional language tasks seen in LLMs, it is intriguing to observe that scaling laws [49, 21, 27], analogous to those observed in LLMs, are also applicable to GAIA-1. This suggests the broader applicability of scaling principles in modern AI models across diverse domains, including autonomous driving.    *

[Hu, Anthony, et al. "GAIA-1: A Generative World Model for Autonomous Driving." *arXiv preprint arXiv:2309.17080* (2023)]

⇒ Idea of the talk:

Let's explore speech and language tech first, namely:

(Section 1) E2E ASR

(Section 2) LLMs

(Section 3) GAIA-1 (finally, knowing transformers already in depth)

* Not wrong, but misleading!
While LLMs have same input and output tokens, GAIA-1 world model doesn't: The input is a multimodal token, thereby asking for an ENC/DEC transformer model (as in E2E ASR)

Technische Universität Braunschweig

Institut für Nachrichtentechnik

Autoregressive decoding of output sequence tokens, token-by-token …

Token, here: ID for letters/digits/signs (~40), words (~300000),
or for so-called byte-pair encodings (subwords) (30000…50000)

Decoded output:
ROCK AND ROLL

The *entire* feature sequence ($\mathbf{x}_1^T$, from $1...T$) is first encoded
into a *hidden* representation sequence $\mathbf{h}_1^T$ of same length
$\Rightarrow$ Streaming not possible, would require modifications

At each decoder token timestep, the decoder uses the
attention function to gather relevant timesteps
from the hidden representation

AED models perform sequence-to-sequence mapping:
$$\mathbf{x}_1^T \to \mathbf{h}_1^T \to y_1, y_2, ..., y_\ell, ..., y_L$$

AED models require large amounts of training data, but
achieve state-of-the-art performance on several datasets

Among the common architectures is the
all-attention-based transformer model
[Vaswani et al., „*Attention is All You Need*", arXiv:1706.03762, 2017]

decoder
token timestep
$\ell =$

hidden =
encoded
representation

output tokens

$y_\ell$

argmax

$\mathbf{h}_1^T$

Encoder

Attention
Decoder

T

$\mathbf{x}_1^T$

$y_{\ell-1}$

frequency

energy

time
feature sequence

previous
output token

Technische
Universität
Braunschweig

Institut für Nachrichtentechnik

No architectural recurrency at all, bypasses used
⇒ very deep models possible

- Encoder consists of encoder blocks using (linear) self-attention and (non-linear) FC layers with residual bypasses

- Decoder in addition uses cross attention (also called: encoder-decoder attention)

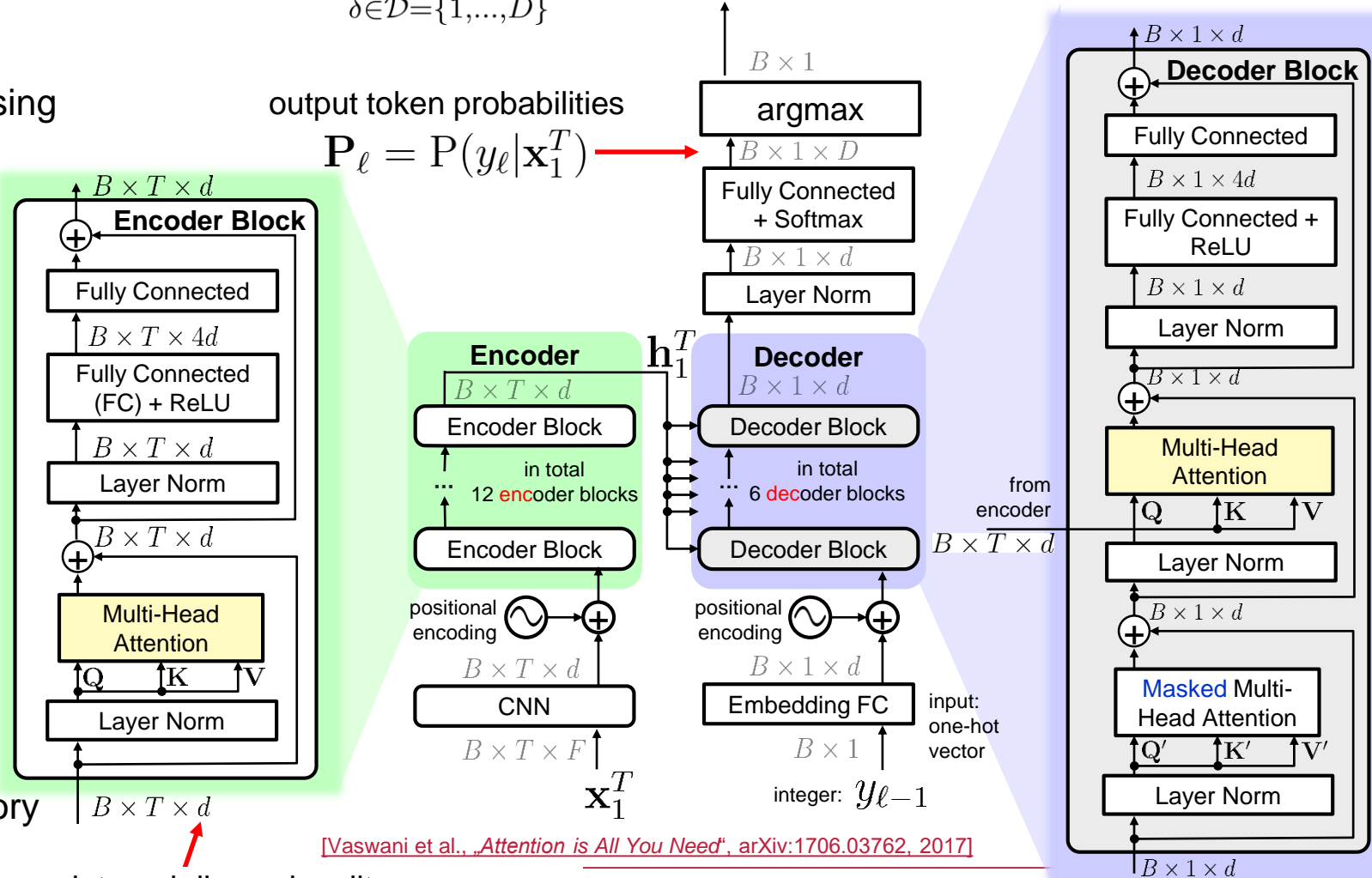- Positional information is lost in the attention layers, therefore, position needs to be encoded both in the encoder and decoder input

- Masked MHA:
  In training, future time steps are masked to zero. In inference, previous timesteps are read from internal memory

$$\arg\max_{\delta \in \mathcal{D} = \{1, \dots, D\}} (\mathrm{P}_{\ell, \delta}) = y_\ell$$

output token probabilities
$$\mathbf{P}_\ell = \mathrm{P}(y_\ell | \mathbf{x}_1^T)$$

$B \times 1$

**argmax**

$B \times 1 \times D$

**Fully Connected + Softmax**

$B \times 1 \times d$

**Layer Norm**

**Encoder** $\mathbf{h}_1^T$   **Decoder**

$B \times T \times d$     $B \times 1 \times d$

Encoder Block     Decoder Block

in total        in total
12 encoder blocks   6 decoder blocks

Encoder Block     Decoder Block     $B \times T \times d$

positional encoding ⊕     positional encoding ⊕     from encoder

$B \times T \times d$     $B \times 1 \times d$

CNN        Embedding FC    input: one-hot vector

$B \times T \times F$     $B \times 1$

$\mathbf{x}_1^T$        integer: $y_{\ell-1}$

**Encoder Block** $B \times T \times d$

⊕

Fully Connected

$B \times T \times 4d$

Fully Connected (FC) + ReLU

$B \times T \times d$

Layer Norm

$B \times T \times d$

⊕

Multi-Head Attention

**Q    K    V**

Layer Norm

$B \times T \times d$

**Decoder Block** $B \times 1 \times d$

⊕

Fully Connected

$B \times 1 \times 4d$

Fully Connected + ReLU

$B \times 1 \times d$

Layer Norm

$B \times 1 \times d$

⊕

Multi-Head Attention

**Q    K    V**

Layer Norm

$B \times 1 \times d$

⊕

Masked Multi-Head Attention

**Q'    K'    V'**

Layer Norm

$B \times 1 \times d$

[Vaswani et al., „*Attention is All You Need*", arXiv:1706.03762, 2017]

internal dimensionality (not dependent on $\dim(\mathbf{x}_t)$ )

## MHA function – Self-attention & cross attention

**Encoder self-attention:**
„Which encoder frame timesteps $t$ (input) relate to which other encoder timesteps $t$ (input) relevantly?"

Encoder-**decoder cross** attention:
„Which encoder frame timesteps $t$ (input) are relevant for the current decoder token timestep $\ell$ (output)?"

**Decoder** masked MHA (**self**-attention):
„Which already decoded token timesteps $1,...,\ell-1$ (previous outputs) are relevant for the current decoder token timestep $\ell$ (output)?"



$d$ : dimension of internal representation
both of input features and of tokens

Encoder-decoder (= cross) attention:
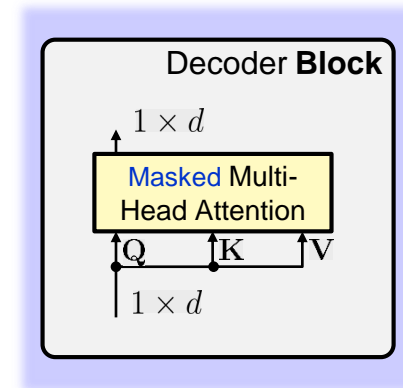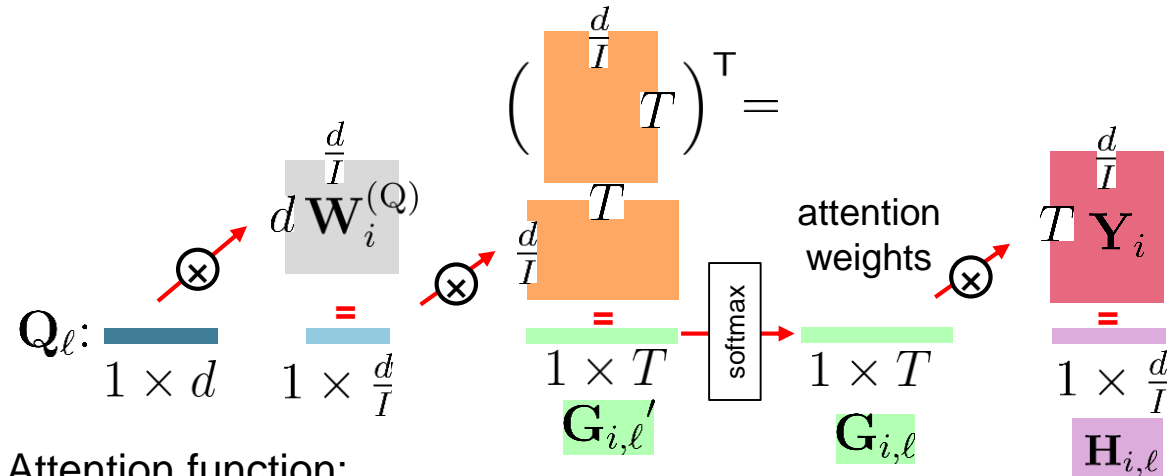„Which encoder frame timesteps $t$ (input) are relevant for the current decoder token timestep $\ell$ (output)?"

MHA employs multiple attention functions in parallel: Attention head $i \in \{1, ..., I\}$

FC layers are simply linear projections

$$d \, \mathbf{W}_i^{(Q)} \quad \frac{d}{I}$$

$$\left( \begin{array}{c} \frac{d}{I} \\ T \end{array} \right)^{\mathsf{T}} =$$

$$\mathbf{Q}_\ell: \quad \underbrace{\quad}_{1 \times d} \quad \otimes \quad \underbrace{\quad}_{1 \times \frac{d}{I}} \quad = \quad \otimes \quad \underbrace{\quad}_{\substack{1 \times T \\ \mathbf{G}_{i,\ell}'}} \xrightarrow{\text{softmax}} \underbrace{\quad}_{\substack{1 \times T \\ \mathbf{G}_{i,\ell}}} \quad \otimes \quad \underbrace{T \, \mathbf{Y}_i}_{\substack{1 \times \frac{d}{I} \\ \mathbf{H}_{i,\ell}}}$$

attention weights

Attention function:

$$\mathbf{H}_{i,\ell} = \mathrm{softmax}\left( \frac{\mathbf{Q}_\ell \mathbf{W}_i^{(Q)} \left( \mathbf{K} \mathbf{W}_i^{(K)} \right)^{\mathsf{T}}}{\sqrt{d}} \right) \mathbf{V} \mathbf{W}_i^{(V)}$$

head index

$\underbrace{\qquad\qquad}_{\text{attention weights } \mathbf{G}_{i,\ell}}$ $\underbrace{\quad}_{\text{value projections } \mathbf{Y}_i}$

Attention weights sum up to one over all $T$ input encoder indices: $\sum_t G_{i,\ell,t} = 1$

Multi-Head Attention

$1 \times d$

Fully Connected

Concatenation $\mathbf{H}_1 | \mathbf{H}_2 | ...$

$\mathbf{H}_{i,\ell}$

Matrix Multiplication

$\mathbf{G}_{i,\ell}$

$\mathbf{Y}_i$

Softmax — Attention Weights

Matrix Multiplication

$1 \times \frac{d}{I} \otimes \frac{1}{\sqrt{d}} \quad \frac{d}{I} \times T \quad ()^{\mathsf{T}} \quad T \times \frac{d}{I}$

FC, $\mathbf{W}_i^{(Q)}$ — FC, $\mathbf{W}_i^{(K)}$ — FC, $\mathbf{W}_i^{(V)}$

from encoder $\quad T \times d$

$\mathbf{K}$ key

$\mathbf{V}$ value

$1 \times d \quad \mathbf{Q}_\ell$ query
from previous decoder block layer

No assumption of a monotonic (= diagonal) input-to-output sequence mapping: Result can be non-monotonic!

Softmax distribution $\mathbf{G}_{i,\ell}$ for decoder index $\ell = 1$

Attended segments in the input speech

Technische
Universität
Braunschweig

Institut für Nachrichtentechnik

$\mathbf{P}_\ell^{\mathrm{ED}}$ ← token probability vectors → $\mathbf{P}_\ell^{\mathrm{LM}}$

$B \times 1 \times D$

Fully Connected + Softmax

$B \times 1 \times d$

Layer Norm

**from encoder**

$\mathbf{h}_{1:T}^{\mathrm{ENC}}$

**Decoder**

Decoder Block

in total ⋯ 6 DEC blocks

*separately trained*

Decoder Block

positional encoding ⊕

Embedding FC

Language Model (LM)

$y_{\ell-1}$

$y_{\ell-1}$

The encoder-decoder ASR transformer model is trained on paired data (audio and text)
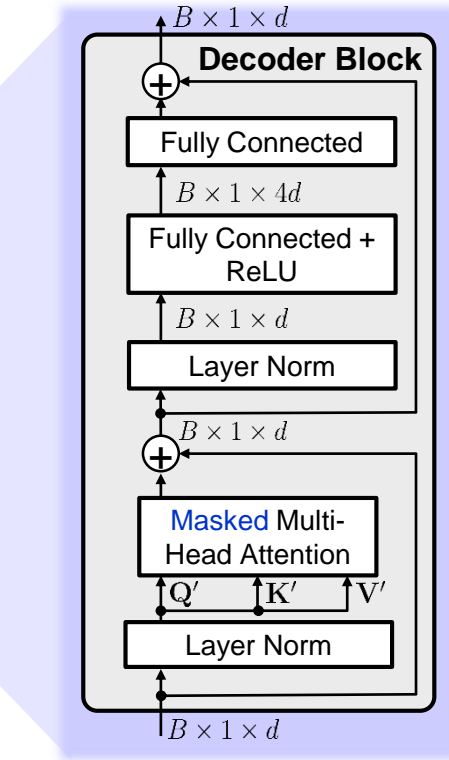
The language model is trained on huge amounts of unpaired data (text only)

Shallow fusion = Use of a LM in the context of E2E ASR:

For each decoding time step $\ell$ :
Combination of output probabilities during inference:

$$\log \mathbf{P}_\ell = \log \mathbf{P}_\ell^{\mathrm{ED}} + \lambda \log \mathbf{P}_\ell^{\mathrm{LM}}$$

~posterior        likelihood        prior

Final output token decision:

$$y_\ell = \underset{\delta \in \mathcal{D} = \{1,...,D\}}{\arg\max} (\mathrm{P}_{\ell,\delta})$$

output token probs

The language model weight $\lambda$ balances the LM contribution

Word error rates with or without external LM:

| | Bidirectional | | Unidirectional | |
|---|---|---|---|---|
| | + | ⊙ | + | ⊙ |
| No external LM | 17.8 | 17.6 | 28.0 | 26.2 |
| Shallow fusion | 15.2 | 13.9 | 22.9 | 21.4 |

From ChatGPT to GAIA-1: On Generative Sequence Models in Speech, Language, and Vision | Prof. Tim Fingscheidt | DEC 6, 2023 | 11

Technische Universität Braunschweig

Institut für Nachrichtentechnik

[G. Gulcehre et al., "*On Using Monolingual Corpora in Neural Machine Translation*", 2015]

[E. McDermott et al., "*A Density Ratio Approach to Language Model Fusion in End-to-End Automatic Speech Recognition*", Proc. of ASRU, 2019]

# 2. Language Model (LM)
## A decoder-only model

From now on, in Section 2, we look at LMs only!
And we look on their usefulness standalone!

$\mathbf{P}_{\ell}^{\mathrm{ED}}$   output token probs

$B \times 1 \times D$

Fully Connected + Softmax

$B \times 1 \times d$

Layer Norm

from **encoder**

$\mathbf{h}_{1:T}^{\mathrm{ENC}}$

**Decoder**

Decoder Block

in total
6 DEC blocks

Decoder Block

positional encoding

Embedding FC

$y_{\ell-1}$

Recap: Original transformer
encoder-decoder ASR model [1]

[1] Vaswani et al., „Attention is All You Need",
arXiv:1706.03762, 2017

---

$\mathbf{P}_{\ell}^{\mathrm{LM}}$

$B \times 1 \times D$

Fully Connected + Softmax

$B \times 1 \times d$

Layer Norm

**Transformer LM**

**Decoder**

Decoder Block

in total
6 DEC blocks

Decoder Block

positional encoding

Embedding FC

$y_{\ell-1}$

Here: Transformer
language model [2]
using decoder blocks

[2] Liu et al., „Generating Wikipedia by Summarizing Long Sequences". ICLR, 2018

---

$B \times 1 \times d$

**Decoder Block**

Fully Connected

$B \times 1 \times 4d$

Fully Connected + ReLU

$B \times 1 \times d$

Layer Norm

$B \times 1 \times d$

Masked Multi-Head Attention

$\mathbf{Q}'$  $\mathbf{K}'$  $\mathbf{V}'$

Layer Norm

$B \times 1 \times d$

Transformer decoder block
in decoder-only transformer
language model [2]

---

Removed:
encoder-decoder
multi-head
cross attention!

Kept:
multi-head
self-attention

We show that generating English Wikipedia articles can be approached as a multi-document summarization of source documents. We use extractive summarization to coarsely identify salient information and a neural abstractive model to generate the article. For the abstractive model, we introduce a decoder-only architecture that can scalably attend to very long sequences, much longer than typical encoder-decoder architectures used in sequence transduction. We show that this model can generate fluent, coherent multi-sentence paragraphs and even whole Wikipedia articles. When given reference documents, we show it can extract relevant factual information as reflected in perplexity, ROUGE scores and human evaluations.

Technische
Universität
Braunschweig

Institut für Nachrichtentechnik

$$y_\ell = \underset{\delta \in \mathcal{D} = \{1,...,D\}}{\arg\max} \left( \mathrm{P}^{\mathrm{LM}}_{\ell,\delta} \right)$$

$\mathbf{P}^{\mathrm{LM}}_{\ell}$    output token probs



$B \times 1 \times D$

Unembedding FC + Softmax

$B \times 1 \times d$

Layer Norm

**Decoder**

Decoder Block

in total
··· **12 DEC blocks**

Decoder Block

$B \times \Lambda \times 1$

positional encoding $\bigoplus$

$B \times \Lambda \times 1$

Embedding FC

$B \times \Lambda \times 1$

$y_{\ell-\Lambda:\ell-1}$

Transformer LM

During pre-training, the previous predictions with a context window of length $\Lambda$ are used to predict the next token probabilities $\mathbf{P}^{\mathrm{LM}}_{\ell}$

Generative pre-trained transformer (GPT) language model [1]:

Unsupervised pre-training by next token prediction:
Loss function: $J = \sum_{\ell \in \mathcal{L}} J_\ell = -\sum_{\ell \in \mathcal{L}} \log \mathrm{P}^{\mathrm{LM}}_{\ell,\delta=\overline{y}_\ell}$ ← token probability of correct token $\overline{y}_\ell$

with the decoding step $\ell \in \mathcal{L} = \{1, 2, ..., L\}$
and $L$ is the length of the ground truth sequence $(\overline{y}_\ell) = \overline{y}_{1:L}$

The GPT-1 [1] language model …
  … has a context window of length $\Lambda = 512$
  … uses byte-pair encoding (BPE), the vocabulary size is 40k
  … consists of 12 decoder blocks with in total 117M parameters
  … requires 1 month on 8 GPUs for pre-training ← A lot! But still university-grade…
  … is pre-trained by publicly available 7000 books
  … code is published on GitHub:
https://github.com/openai/finetune-transformer-lm

[1] Radford et al. "Improving Language Understanding by Generative Pre-Training." (2018).

Technische Universität Braunschweig

Institut für Nachrichtentechnik

## The rise of GPT: From GPT-1 to GPT4

Only applied in discriminative tasks (i.e., classification)

Also used for generative tasks, e.g., writing stories, neural machine translation

Good performance in zero-shot and few-shot settings

Multi-modal inputs including image and texts

Multi-lingual processing

**GPT-1 [1]: 12 dec blocks**
- Context window length: $\Lambda = 512$
- Training data: BookCorpus with 7000 books
- #params: 117 million = $1.17 \times 10^8$

**GPT-2 [2]: 48 dec blocks**
- Context window length: $\Lambda = 1024$
- Training data: BookCorpus and WebText (8M webpages)
- #params: 1.5 billion = $1.5 \times 10^9$

**GPT-3 [3]: 96 dec blocks**
- Context window length: $\Lambda = 2048$
- Training data: Common Crawl (410B tokens) and WebText2 (19B tokens)
- #params: 175 billion = $1.75 \times 10^{11}$

**GPT-4 [4]: ?? dec blocks**
- Context window length: $\Lambda = 32768$
- Training data: not written in [4]
- #params: ~100 trillion = $1 \times 10^{14}$

Generative pre-training …

… + pre-training with task conditioning …

… + pre-training as GPT-2 + in-context learning

[1] Radford et al. "Improving Language Understanding By Generative Pre-training." (2018).
[2] Radford et al. "Language Models Are Unsupervised Multitask Learners." *OpenAI blog* 1.8 (2019): 9.
[3] Brown et al. "Language Models Are Few-Shot Learners." *in Proc. of NeurIPS*, virtual, Dec, 2020, 1877-1901.
[4] OpenAI "GPT-4 Technical Report." *arXiv* (2023): 2303-08774.

# GPT-4: Multi-modal large language model
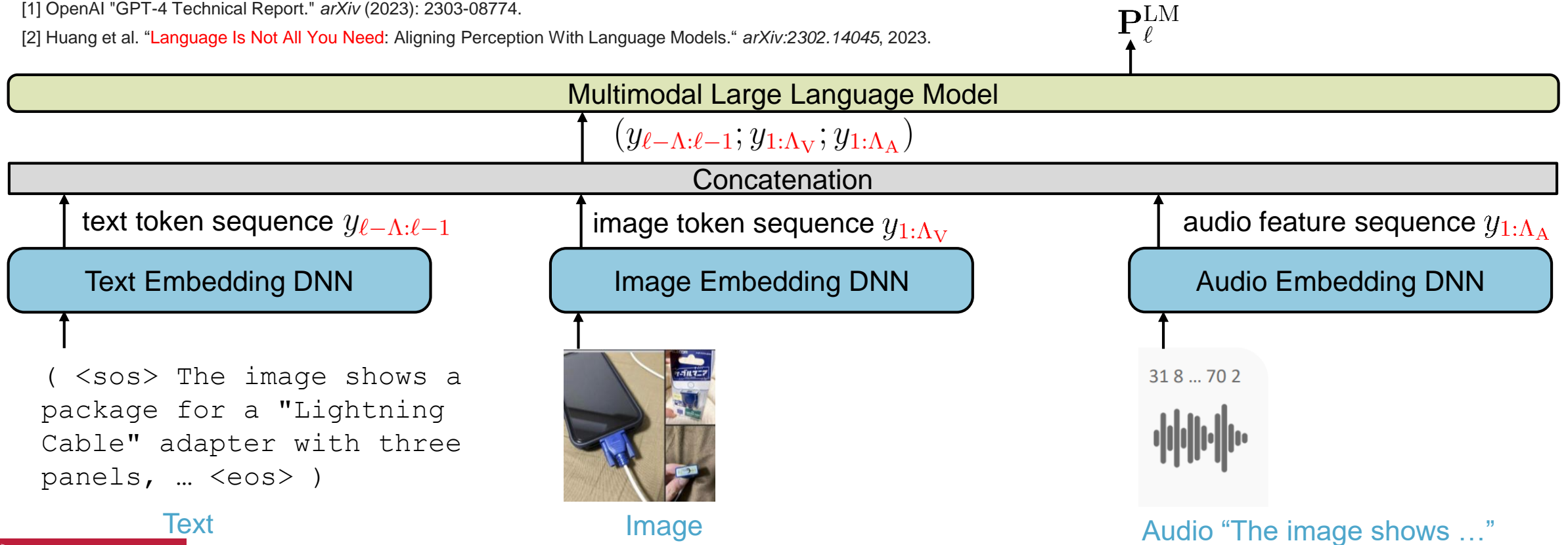
GPT4 can take both text and image as input.

The technical report from OpenAI [1] doesn't give any details on model architecture and training.

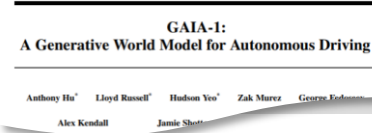The following technique is based on a different multi-modal large language model (LLM) [2]:

[1] OpenAI "GPT-4 Technical Report." *arXiv* (2023): 2303-08774.

[2] Huang et al. "Language Is Not All You Need: Aligning Perception With Language Models." *arXiv:2302.14045*, 2023.

$\mathbf{P}_{\ell}^{\mathrm{LM}}$

Multimodal Large Language Model

$$\left( y_{\ell-\Lambda:\ell-1}; y_{1:\Lambda_{\mathrm{V}}}; y_{1:\Lambda_{\mathrm{A}}} \right)$$

Concatenation

text token sequence $y_{\ell-\Lambda:\ell-1}$  image token sequence $y_{1:\Lambda_{\mathrm{V}}}$  audio feature sequence $y_{1:\Lambda_{\mathrm{A}}}$

Text Embedding DNN  Image Embedding DNN  Audio Embedding DNN

```
( <sos> The image shows a
package for a "Lightning
Cable" adapter with three
panels, … <eos> )
```



31 8 … 70 2

Text        Image        Audio "The image shows …"

The blue arrow means the modality can be dropped in inference

**Tokenizer**

quantized features

latent representations:

**Image Tokenizer**
trained
2D U-Net ENC [31]
and quantization [28]
#params: 0.3B

$\mathbf{x}_1^T = (\mathbf{x}_t)$

Video

$\mathbf{x}_t \in \mathbb{R}^{H \times W \times 3}$

$\mathbf{z}_1^T$

$\mathbf{z}_t \in \mathbb{R}^{n \times d}$

$n = \frac{H}{16} \times \frac{W}{16} = 576$ tokens/frame

FC($d$)
jointly trained with world model

$\tilde{\mathbf{a}}_1^T = (\tilde{\mathbf{a}}_t)$

Action
(speed, curvature)

$\tilde{\mathbf{a}}_t \in \mathbb{R}^{2 \times 1}$

$\mathbf{a}_1^T$

$\mathbf{a}_t \in \mathbb{R}^{2 \times d}$

2 tokens/frame

**Text Encoder**
pretrained/fixed
transformer ENC/DEC
#params: 0.77B [24]

$\tilde{\mathbf{c}}_1^T = (\tilde{\mathbf{c}}_{t,m})$

Text

$\tilde{\mathbf{c}}_{t,m} \in \mathcal{C}$

$\mathbf{c}_1^T$

$\mathbf{c}_t \in \mathbb{R}^{M \times d}$

$M = 32$ tokens/frame

**Conditions**

token set

**Concat & Positional Embedding**

$\mathbf{b}_1^T$

$\mathbf{b}_t = \text{pos}(\text{concat}(\mathbf{c}_t, \mathbf{z}_t, \mathbf{a}_t))$
$\mathbf{b}_t \in \mathbb{R}^{(M+n+2) \times d} = \mathbb{R}^{610 \times d}$

temporally predictive image embeddings

**World Model**
trained
transformer ENC/DEC
#params: 6.5B

$\hat{\mathbf{z}}_{T+1}^{T+N}$

$N$: # of future frames to be predicted

frame rate = 6.25 fps (index $t$)

frame rate = 25 fps (index $\tau$)

Interpolation capability: Between any 2 frames, fill in 3 interpolated frames

**Video Model**
trained
3D U-Net [38]
ENC/DEC
video diffusion model
#params: 2.6B

$\hat{\mathbf{x}}_{4T+4}^{4T+4N}$

$\hat{\mathbf{x}}_\tau \in \mathbb{R}^{H \times W \times 3}$

Recover full image resolution (x16 rows and columns)

[24] Raffel et al., "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer". Journal of Machine Learning Research, 2020. Here: **T5-large** model.
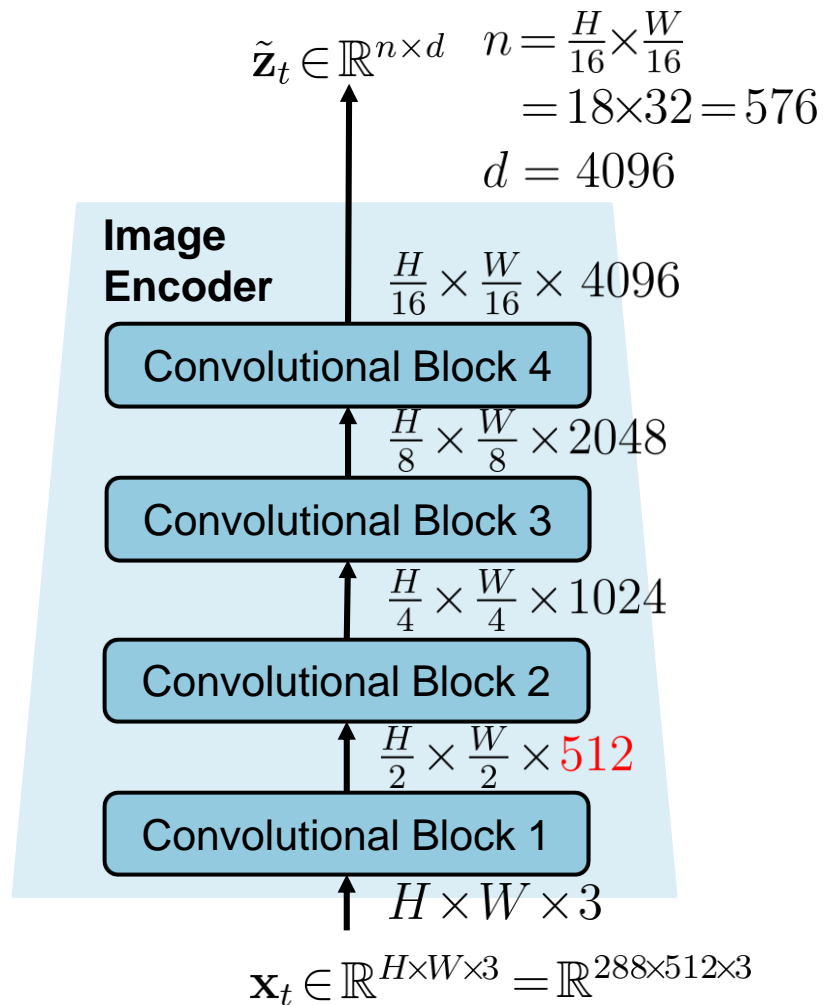[28] Oord et al., „Neural Discrete Representation Learning". In Proc. of NeurIPS, 2017.
[31] Ronneberger et al. „U-Net: Convolutional Networks for Biomedical Image Segmentation". In Medical Image Computing and Computer-Assisted Intervention (MICCAI), 2015.
[38] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video Diffusion Models." arXiv, Jun. 22, 2022.

Technische Universität Braunschweig

Institut für Nachrichtentechnik

Image tokenizer (training)



**Image Tokenizer** (#params: 0.3B)

quantizer codebook, learned, with low dimension $d' < d$

$32 < 4096$ expanded codevector dimension

$\mathbb{D}' = \mathbb{R}^{K \times d'}$

# of CB entries: $K = 8192$

$\mathbf{x}_t \in \mathbb{R}^{H \times W \times 3} = \mathbb{R}^{288 \times 512 \times 3}$

Image

$\mathbf{x}_t$

Image Encoder ~[31]

$\tilde{\mathbf{z}}_t \in \mathbb{R}^{n \times d}$

Linear FC VQ, Linear FC [28]

$\mathbb{D} = \mathbb{R}^{K \times d}$

$\mathbf{z}_t \in \mathbb{D}^n \subset \mathbb{R}^{n \times d}$

Image Decoder (only used for training)

$\hat{\hat{\mathbf{x}}}_t$

$J^{\mathrm{rec}}$ (1)

target: $\mathbf{x}_t \in \mathbb{R}^{H \times W \times 3}$

Reconstruction loss: L1, L2, perceptual loss, and GAN loss

$n = \frac{H}{16} \times \frac{W}{16} = 18 \times 32 = 576$

target

$J^{\mathrm{quant}}$ (2)

reconstructs the image $\hat{\hat{\mathbf{x}}}_t$ from quantized image embeddings $\mathbf{z}_t$

Quantization loss (embedding loss): commitment loss and L2 loss

$J^{\mathrm{IB}}$ (3)

Inductive bias loss: cosine similarity loss

Quantized latent representation shall be similar to SSL DINO representations

DINO SSL model [30] pretrained vision model #params: **21/23/85M**

target

**frozen**

The final loss has 3 components: (1) image reconstruction loss, (2) quantization loss, (3) inductive bias loss

[28] Oord et al., „Neural Discrete Representation Learning". In Proc. of NeurIPS, 2017.

[31] Ronneberger et al. „U-Net: Convolutional Networks for Biomedical Image Segmentation". In Medical Image Computing and Computer-Assisted Intervention (MICCAI), 2015.

[30] Caron et al., „Emerging Properties in Self-Supervised Vision Transformers". In Proc. of ICCV, 2021.
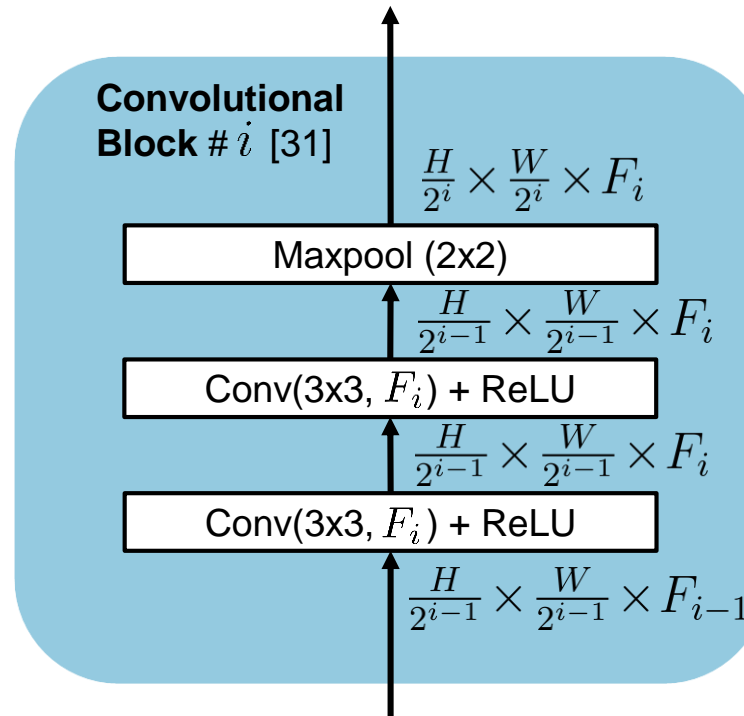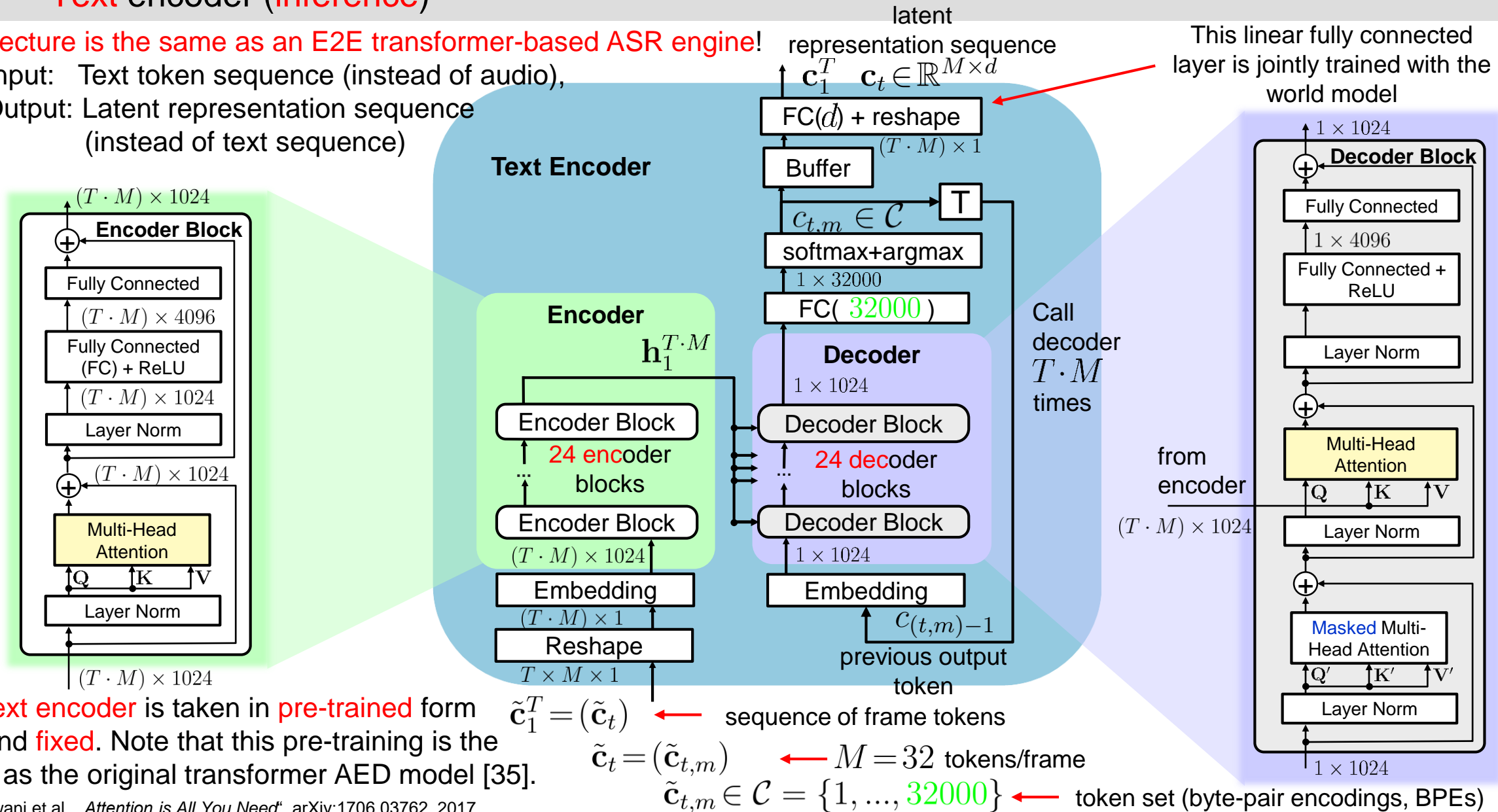
Technische Universität Braunschweig

Institut für Nachrichtentechnik

## Image encoder

$\tilde{\mathbf{z}}_t \in \mathbb{R}^{n \times d}$   $n = \frac{H}{16} \times \frac{W}{16}$
$= 18 \times 32 = 576$
$d = 4096$

**Image Encoder**

$\frac{H}{16} \times \frac{W}{16} \times 4096$

Convolutional Block 4

$\frac{H}{8} \times \frac{W}{8} \times 2048$

Convolutional Block 3

$\frac{H}{4} \times \frac{W}{4} \times 1024$

Convolutional Block 2

$\frac{H}{2} \times \frac{W}{2} \times 512$

Convolutional Block 1

$H \times W \times 3$

$\mathbf{x}_t \in \mathbb{R}^{H \times W \times 3} = \mathbb{R}^{288 \times 512 \times 3}$

*„The discrete autoencoder is a fully convolutional U-Net structure [31]"*

- However, no architecture details to the image encoder in GAIA-1
- Here: Reverse engineering: The output dimension of the 1st convolutional block can be changed from 64 to 512 to match the 0.3B parameters written in GAIA-1

**Convolutional Block #$i$ [31]**

$\frac{H}{2^i} \times \frac{W}{2^i} \times F_i$

Maxpool (2x2)

$\frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}} \times F_i$

Conv(3x3, $F_i$) + ReLU

$\frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}} \times F_i$

Conv(3x3, $F_i$) + ReLU

$\frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}} \times F_{i-1}$

In each convolutional block [31] do:
- Upsample the feature dimension by 2 (except the 1st one) $F_i = 2F_{i-1}$
- Downsample each spatial dimension by 2

[31] Ronneberger et al. „U-Net: Convolutional Networks for Biomedical Image Segmentation". In Medical Image Computing and Computer-Assisted Intervention (MICCAI), 2015.

Technische Universität Braunschweig

Institut für Nachrichtentechnik

## Text encoder (inference)

Architecture is the same as an E2E transformer-based ASR engine!
But: Input:   Text token sequence (instead of audio),
     Output: Latent representation sequence
             (instead of text sequence)

This linear fully connected layer is jointly trained with the world model



latent representation sequence
$\mathbf{c}_1^T \quad \mathbf{c}_t \in \mathbb{R}^{M \times d}$

FC($d$) + reshape

$(T \cdot M) \times 1$

Buffer

T

$c_{t,m} \in \mathcal{C}$

softmax+argmax

$1 \times 32000$

FC( 32000 )

Call decoder $T \cdot M$ times

**Text Encoder**

**Encoder**

$\mathbf{h}_1^{T \cdot M}$

**Decoder**

$1 \times 1024$

Encoder Block

24 encoder blocks

Decoder Block

24 decoder blocks

Encoder Block

$(T \cdot M) \times 1024$

Decoder Block

$1 \times 1024$

Embedding

$(T \cdot M) \times 1$

Embedding

$C_{(t,m)-1}$

Reshape

$T \times M \times 1$

previous output token

**Encoder Block** $(T \cdot M) \times 1024$

Fully Connected

$(T \cdot M) \times 4096$

Fully Connected (FC) + ReLU

$(T \cdot M) \times 1024$

Layer Norm

$(T \cdot M) \times 1024$

Multi-Head Attention

Q  K  V

Layer Norm

$(T \cdot M) \times 1024$

**Decoder Block**

$1 \times 1024$

Fully Connected

$1 \times 4096$

Fully Connected + ReLU

Layer Norm

Multi-Head Attention

from encoder

Q  K  V

$(T \cdot M) \times 1024$

Layer Norm

Masked Multi-Head Attention

Q′  K′  V′

Layer Norm

$1 \times 1024$

The text encoder is taken in pre-trained form [24] and fixed. Note that this pre-training is the same as the original transformer AED model [35].

$\tilde{\mathbf{c}}_1^T = (\tilde{\mathbf{c}}_t)$ ← sequence of frame tokens

$\tilde{\mathbf{c}}_t = (\tilde{\mathbf{c}}_{t,m})$ ← $M = 32$ tokens/frame

$\tilde{\mathbf{c}}_{t,m} \in \mathcal{C} = \{1, ..., 32000\}$ ← token set (byte-pair encodings, BPEs)

[35] Vaswani et al., „*Attention is All You Need*", arXiv:1706.03762, 2017

[24] Raffel et al., "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer". Journal of Machine Learning Research, 2020. Here: **T5-large** model.

Concatenation & positional embedding

Architecture is again similar to an E2E transformer-based ASR engine!

1 output modality: image (tokens) $\hat{\mathbf{z}}_{T+1}^{T+N}, \quad \hat{\mathbf{z}}_t = (\hat{\mathbf{z}}_{t,\nu}) \in \mathbb{D}^n \subset \mathbb{R}^{n \times d}$ temporally predictive image embeddings

**World Model**

$N \times n \times d$

CB Lookup+Reshape

$(N \cdot n) \times 1$

Buffer

1

argmax

$\mathbf{P}_{t,\nu} \in [0,1]^K$

$K = 8192$

softmax

"negative text prompting"

$K = 8192$

$-\alpha_{t,\nu}$ ⊗ w/o text $\mathbf{c}_t$ prompting

$t \in \{T+1, ..., T+N\}$
# of future frames to be predicted

⊕

$1 + \alpha_{t,\nu}$ ⊗

T

$\nu \in \{1, ..., n\}$
# of image tokens per frame

**Encoder**

$\mathbf{g}_1^{T \cdot 610}$

FC(8192)

$N \cdot n$ times AR decoder call

Encoder Block

$1 \times d$ **Decoder**

Decoder Block

from encoder

Encoder Block

··· encoder blocks

··· decoder blocks

$(T \cdot 610) \times d$

$(T \cdot 610) \times d$

Decoder Block

$\mathbf{P}_{(t,\nu)-1}$

Encoder Block

$1 \times d$

$(T \cdot 610) \times d$

$\hat{\mathbf{z}}_{(t,\nu)-1} \in \mathbb{D}$

Random Top-k Sampling

Reshape

CB Lookup

$T \times 610 \times d$

$\hat{k} \in \{1, 2, ..., 8192\}$

The training process of the world model is the same as the original transformer encoder-decoder model [35]

$(T \cdot 610) \times d$

Encoder Block

Fully Connected

$(T \cdot 610) \times 4d$

Fully Connected (FC) + ReLU

$(T \cdot 610) \times d$

Layer Norm

$(T \cdot 610) \times d$

⊕

Multi-Head Attention

Q  K  V

Layer Norm

$(T \cdot 610) \times d$

**Decoder Block**

$1 \times d$

⊕

Fully Connected

$1 \times 4d$

Fully Connected + ReLU

Layer Norm

⊕

Multi-Head Attention

Q  K  V

Layer Norm

$(T \cdot 610) \times d$

⊕

Masked Multi-Head Attention

Q'  K'  V'

Layer Norm

$1 \times d$

$\mathbf{b}_1^T$

previous output image embedding index

[35] Vaswani et al., „Attention is All You Need", arXiv:1706.03762, 2017

The video generation model (3D U–Net) can be conditioned on image embeddings ( $\hat{\mathbf{z}}$ ) and/or on images ($\underline{\hat{\mathbf{x}}}$ , $\underline{\underline{\hat{\mathbf{x}}}}$).

When conditioned on image embeddings,
it serves as a frame generator

When conditioned on images only,
it serves as an interpolator

Operation in chunks of
length $T' = 7$:

frame rate
= 6.25 fps

frame rate
= 12.5 fps

frame rate
= 25 fps

**Video Model**

(no input)

$\underline{\hat{\mathbf{x}}}_{T+6}^{T+7}$
(last 2 frames)

Buffer

entire sequence of
temporally predicted
image embeddings
(from world model)

$\hat{\mathbf{z}}_{T+1}^{T+N}$

Buffer

**3D U-Net** [38]

(as generator)

$\hat{\mathbf{z}}_{T+1}^{T+7}$

$\hat{\mathbf{z}}_{T+8}^{T+12}$

$\underline{\hat{\mathbf{x}}}_{T+1}^{T+7}$

$\underline{\hat{\mathbf{x}}}_{T+7+1}^{T+7+5}$

Buffer

$\underline{\hat{\mathbf{x}}}_{T+1}^{T+N}$

$(\mathbf{0})$

**3D U-Net** [38]

(as interpolator)

$\underline{\underline{\hat{\mathbf{x}}}}_{2T+2}^{2T+2N}$

$(\mathbf{0})$

**3D U-Net** [38]

(as interpolator)

$\hat{\mathbf{x}}_{4T+4}^{4T+4N}$

… until $N$ video frames are generated

first chunk: video generation

second chunk: AR video generation

no embeddings $\hat{\mathbf{z}}$ are used as condition

[38] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video Diffusion Models." arXiv, Jun. 22, 2022.

Technische Universität Braunschweig

From ChatGPT to GAIA-1: On Generative Sequence Models in Speech, Language, and Vision | Prof. Tim Fingscheidt | DEC 6, 2023 | 22

Institut für Nachrichtentechnik

In training, the video model is conditioned on a sequence of images and image embeddings encoded by the image tokenizer:



A single 3D U-Net model is trained on multiple tasks, conditioned on (masked) image embeddings and (masked) images
The selector positions are specified by the training task:          (each task is equally represented in training)

**(a) Image generation**
(temp. layers deactivated)

**(b) Video generation**

**(c) Autoregressive video generation**
(2 old images and 5 new world model outputs)

**(d) Video interpolation**
(every other image and all image tokens are masked)

[38] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion models. In arXiv preprint, 2022.

frame rate = 6.25 fps

frame rate = 12.5 fps or 25 fps

The neural network architecture is a 3D U-Net with factorized spatial and temporal attention layers [38]

Image/video generation by iterative denoising in 50 steps $i = 1, 2, ..., 50$

$M_j$ defines the feature dimensionality after a 3D U-Net block and can be configured

[38] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video Diffusion Models." arXiv, Jun. 22, 2022.

input to the next
spatial attention block

**Convolutional Residual Block**

$T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$

$+$

$T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$    $T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$

MLP

$T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$

FC( $M_j$ )    BatchNorm + ReLU

$1 \times 1 \times 1 \times M_j$

$T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times d$

MLP    Upscaling / Downscaling    Conv (1x3x3, $M_j$)

$1 \times 1 \times 1 \times M_j$    $T' \times \frac{H}{16} \times \frac{W}{16} \times d$

Reshape    BatchNorm+ ReLU

$M_j$

FC( $M_j$ )    Reshape    Conv (1x3x3, $M_j$)

$1$    $T' \times (\frac{H}{16} \cdot \frac{W}{16}) \times d$    $T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$

$i \in \{1, 2, ..., 50\}$    $\hat{\mathbf{z}}$    output of
diffusion    image    previous block
iteration step    embedding
condition    condition

$n = \frac{H}{16} \times \frac{W}{16} = 576$

**3D U-Net Block**

Upscaling / Downscaling

Temporal    $T' \times H \times W \times M_j$
Attention
Block    Reshape

Multi-Head
Attention

$(H \cdot W) \times T' \times M_j$

Reshape

Spatial
Attention    Multi-Head
Block    Attention

$T' \times (H \cdot W) \times M_j$

Reshape

$i, \hat{\mathbf{z}}$    Convolutional Residual Block

$T' \times H \times W \times M_j$

**Attention Block**

$T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$

Conv (1x1x1)

Matrix Multiply

Softmax

* rel. positional
embedding **B**    $\oplus$

Matrix Multiply

**Q**    **K**    **V**
Conv (1x1x1)    Conv (1x1x1)    Conv (1x1x1)

$T' \times \frac{H}{2^j} \times \frac{W}{2^j} \times M_j$

* only used in
temporal
attention block

[38] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video Diffusion Models." arXiv, Jun. 22, 2022.

Technische
Universität
Braunschweig

Institut für Nachrichtentechnik

Driving data has been collected in London, UK (2019-2023)

<span style="color:red">Training data</span>: 4.700 hours at 25Hz (~420 million images)

Data balancing to account for geography and visually distinct weather conditions:

- Tokenizer balanced over latitude, longitude, weather conditions
- World model & video model balanced over latitude, longitude, weather, steering behaviour, speed behaviour

<span style="color:red">Validation data</span>: 400 hours

Validation within strict predetermined geofences:

- 2 geofences with roads never seen during training
- 1 geofence around the main data collection routes but with runs not used during training

Driving data road map (London)



Train          Validation

## Dataset Sizes

How does GAIA-1 training data compare to typical open access datasets in automotive research?

| Name | # images | # annotated images | total length of videos [h] | # frames per second |
|------|---------|--------------------|---------------------------|---------------------|
| CamVid | 701 | 701 | <1 | 1.00 - 15.00 |
| KITTI | 19,103 | 19,103 | <1 | 10.00 |
| Cityscapes | 150,000 | 5,000 | 5 | 16.67 |
| Waymo Open Perception | 230,000 | 230,000 | 7 | 10.00 |
| A2D2 | 392,556 | 41,277 | <1 | 30.00 |
| Caltech Pedestrian | 1,000,000 | 250,000 | 10 | 30.00 |
| nuScenes | 1,200,000 | 40,000 | 15 | 11.67 |
| SODA10M | 10,000,000 | 20,000 | 27,833 | 0.10 |
| BDD100K | 120,000,000 | 100,000 | 1,111 | 30.00 |
| **GAIA-1 train+val** | **420,000,000** | **?** | **4,700 + 400** | **25.00** |

← BDD is useful!
Annotations for video attributes, include weather, scene, and time of day.

Technische Universität Braunschweig

Institut für Nachrichtentechnik

**3.    GAIA-1**
Further Results …

Using text or (here:) actions, GAIA-1 can be forced to drive onto the pavement (never seen in training!)

GAIA-1 knows: Green light means go ahead!

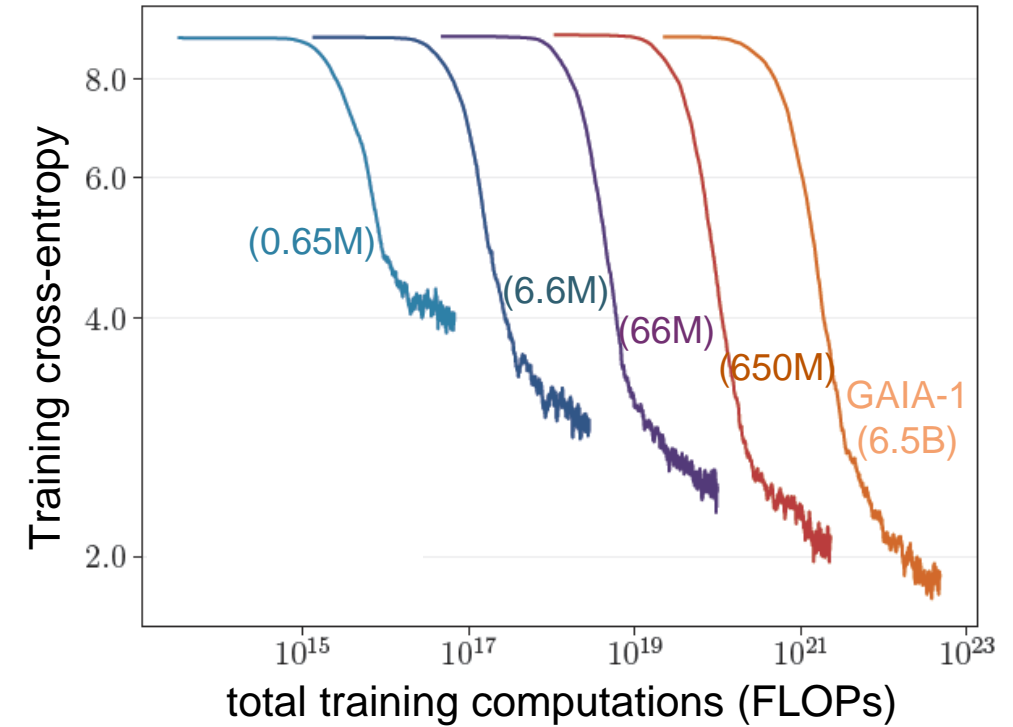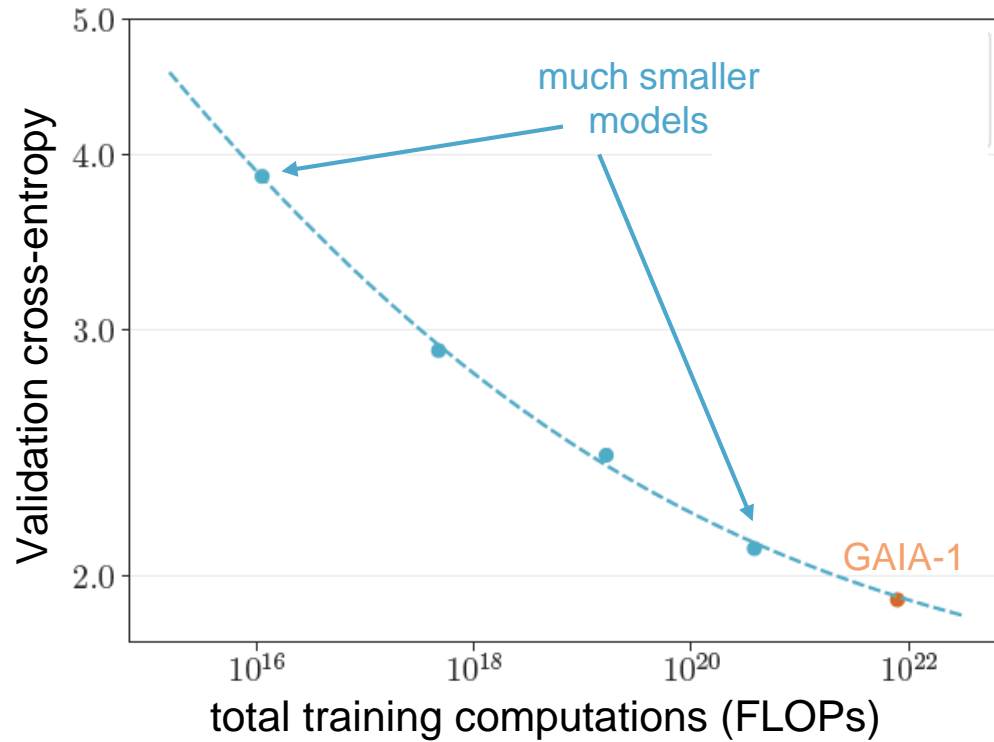No video prompt, but text conditioning (here: weather, daytime)

From ChatGPT to GAIA-1: On Generative Sequen

GAIA-1 world model validation performance is predictable from smaller world models

GAIA-1 world model training performance gets better and better with a larger world model and the use of more data

# Conclusions

ASR: End2end automatic speech recognition achieves SOTA performance with attention-based encoder-decoder (AED) models

LLMs: (Large) language models (e.g., ChatGPT) achieve SOTA with attention-based decoder models

GAIA-1 achieves impressive results with an attention-based encoder-decoder (AED) world model

What we can learn:

Use standard separately trained tokenizers for each input modality; discretize patches of input images

Build multimodal foundational world models, integrating language and vision

Let the world model do the temporal prediction, and …

… let the video model reconstruct the output video in chunks

„Attention is all you need": It seems to be somewhat true…

[Vaswani et al., „*Attention is All You Need*", arXiv:1706.03762, 2017]

Technische
Universität
Braunschweig

Institut für Nachrichtentechnik

# Thank you
# for your attention …

Prof. Dr.-Ing. Tim Fingscheidt

t.fingscheidt@tu-bs.de

Slides on ResearchGate

Technische
Universität
Braunschweig

Institut für Nachrichtentechnik