

# Social Attention Is All We Need

By Ahmad Khalidi

## Abstract

- Do data subjects have responsibilities about how their personal data are used in AI?
- Latest research in responsible AI is not asking that question. Primary responsibility lies within experts processing personal data [1].
- This work proposes an approach, where data subjects map their data on latent spaces..
- These latent spaces are regulated from data subjects personal subjective viewpoints (VP).
- The first approach uses variational encoder decoder transformer in a reinforcement setting.
- This study will try to prove that
  - Unregulated VPs are individual, noisy, entangled, and thus difficult to interpret,
  - Increase global and decrease local expense and
  - help developing responsible AI on large scale.

## Methods

- Drawing from my personal experience and observations, Natural Language Processing (NLP) allows nuanced information exchange, thus motivating the adaption for this study.
- Transformers have shown great results in NLP and have been successfully trained with (self-) supervised and reinforcement learning [2, 3].
- The encoder and decoder models play the referential game [4] in a multi agent cooperative Environment on the MARS service[5].
- The encoders learn to specifically address and exclude selected VPs of the decoders.
- Each training phase requires less reliance on the environment and participating agents.

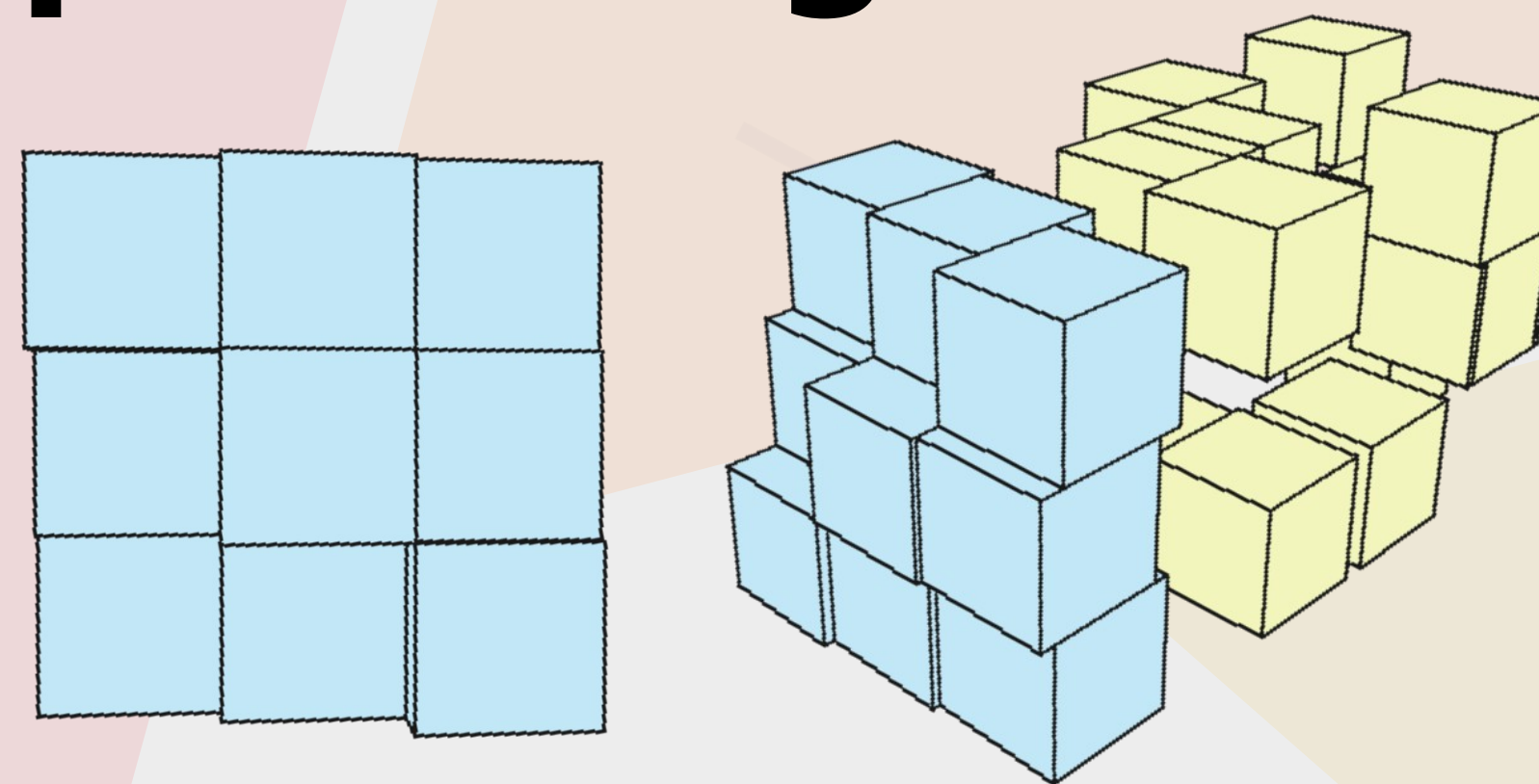
## Results

- Passing forward in BERT[6] and producing a covariance matrix is expensive.
- 8800 last hidden state tensors with labels and logits for four VPs (identity, add, mul, distribution) have been collected for offline training.
- Dynamic architectures in final dense layers are necessary to enable expensive policy.

## Further Research

- Find trade-off between accuracy and efficiency in policy architecture. Potential policy architectures are CNNs with residual blocks and multi head attention layer with down sampling.
- Experiment with phase 2 using PPO and clipped surrogate objective function.
- Research and apply evaluation metrics for the learned privacy enhancing techniques.

# Natural language learned by variational encoder-decoder transformer with viewpoint regularization



## for social driven responsible AI

### Phase 1 Pre-Training

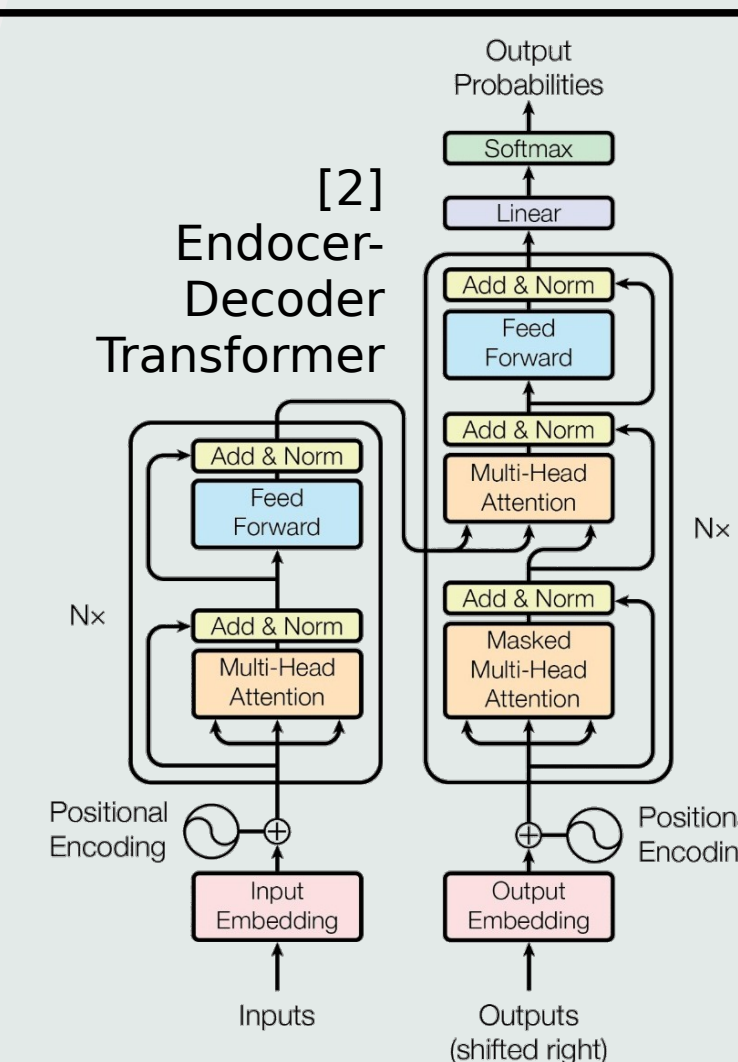
- Single Multi-Head Encoder-Decoder
- One or more heads represent a VP
- Training datapoints:
  - target VP
  - input params
- Decoder and Headers (VPs) are trained to fit datapoints by disregarding target VP
- Policy is trained to behave like the input and a normal Gaussian distribution.

### Phase 2 Encoder Policy

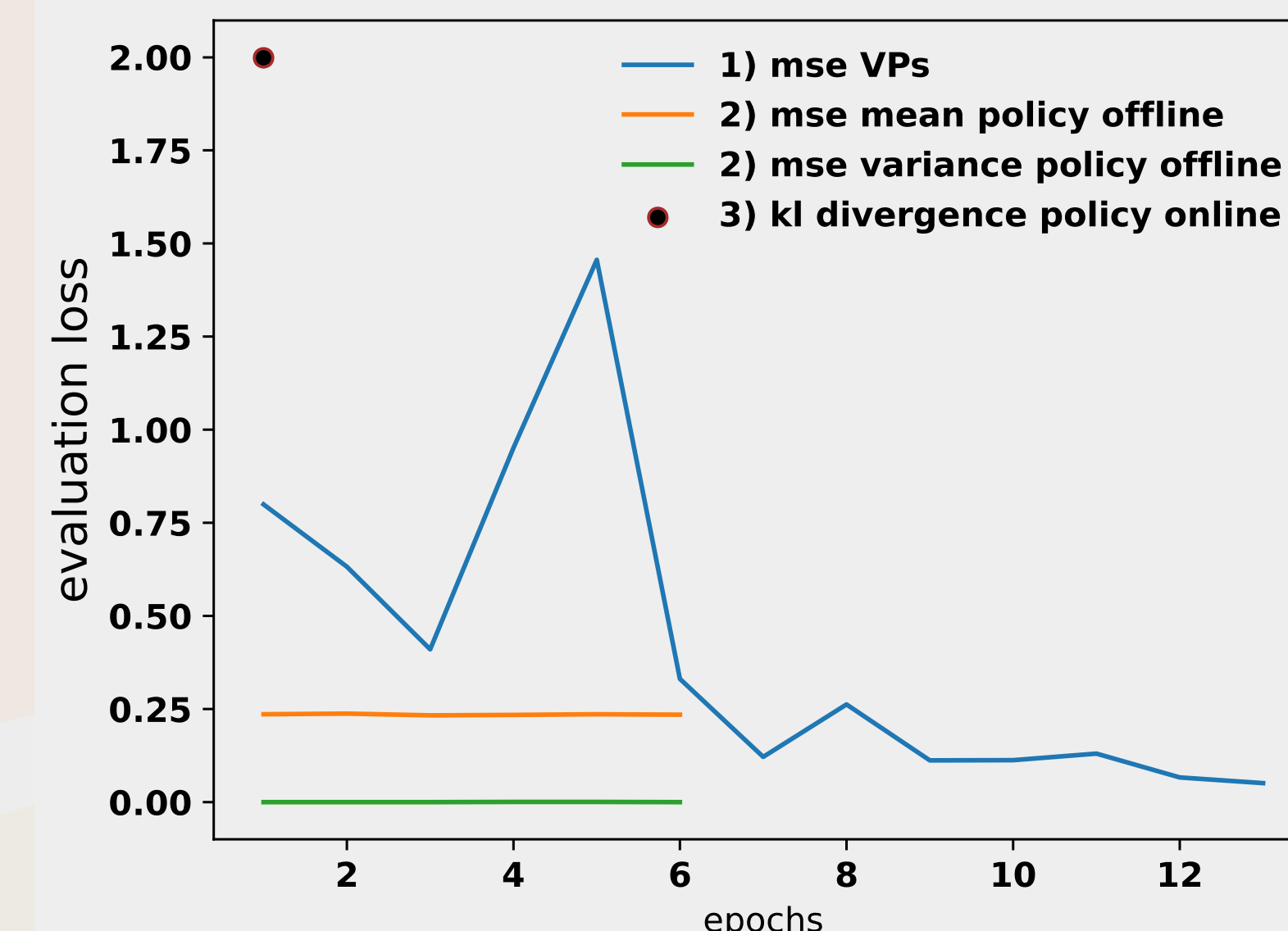
- Separate encoder from decoder
- Header for Encoder:
  - (Co-)Variances
  - Mean
- Sample latent space
- Reward of encoder is a weighted sum of the decoders losses
  - If target VP → negative margin loss
  - If not target VP → positive margin loss

### Phase 3 Decoder Policy

- Similar architecture to Encoder Policy
- Decoders learn what targets and strategy they use to interpret latent space samples

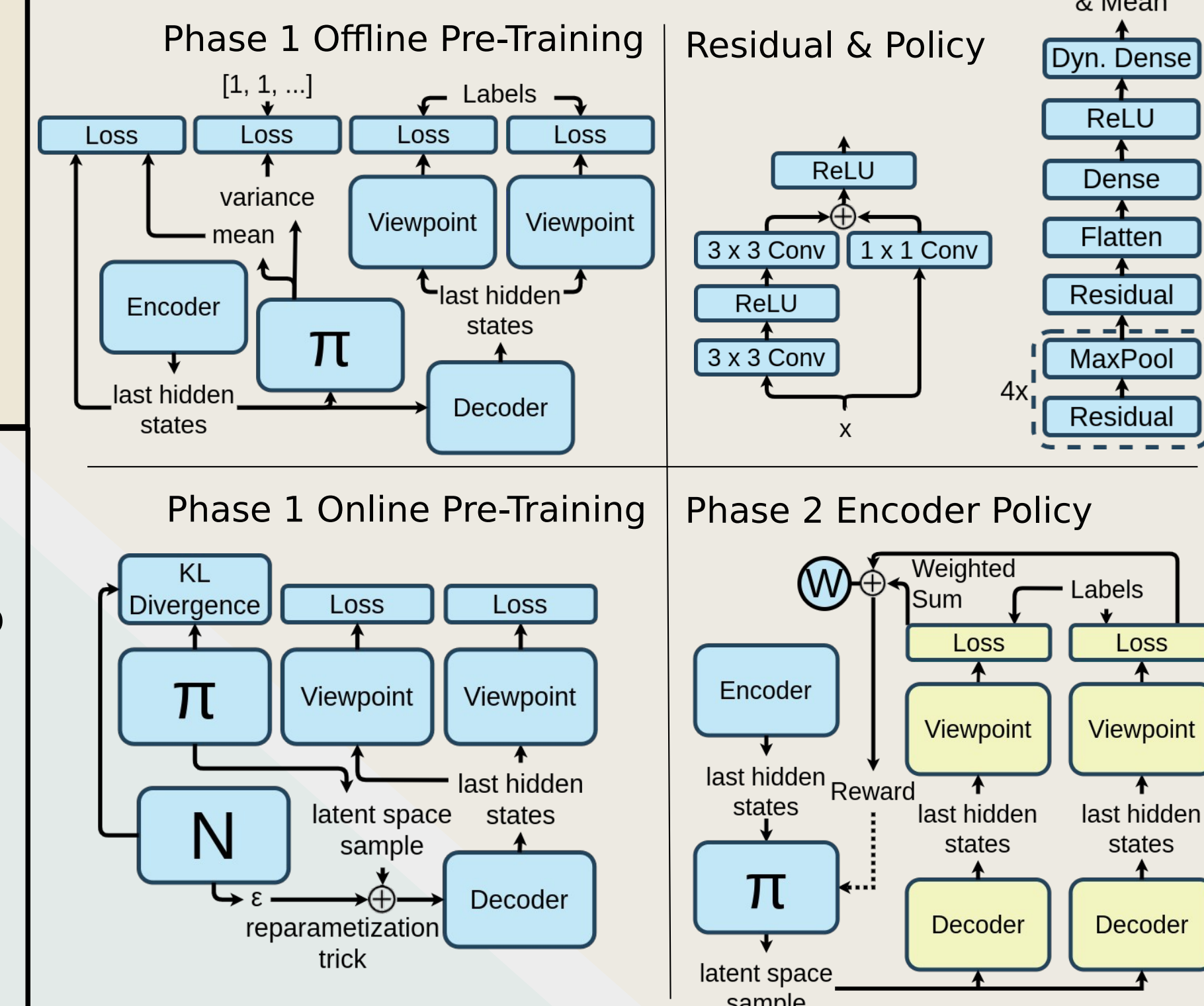


## Phase 1 Pre-Training Progress



• VP's mse policy online: 200.0

## Architectures



## Bibliography

- [1] Göllner, Sabrina et al. „Responsible Artificial Intelligence:Structured Literature Review“ (2023)
- [2] Vaswani, Ashish et al. “Attention is All you Need“ (2017).
- [3] Ziegler, Daniel M. et al. “Fine-Tuning Language Models from Human Preferences.” (2019)
- [4] Lewis, David. “Convention: A Philosophical Study.” (1970).
- [5] Hüning, Christian et al. „Modeling & Simulation as a Service with the Massive Multi-Agent System MARS“ (2016)
- [6] Devlin, Jacob,et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.”(2019)



Visit me on github  
for latest research  
results!

