

ULTra: Unveiling Latent Token Interpretability in Transformer-Based Understanding

Hesam Hosseini, Ghazal Hosseini Mighan, Amirabbas Afzali

{hesam8hosseini, ghazaldesu, amir8afzali}@gmail.com

Sajjad Amini,

samini@umass.edu, s.amini@sharif.edu

Amir Houmansadr

amir@cs.umass.edu

Abstract

Transformers have revolutionized Computer Vision (CV) and Natural Language Processing (NLP) through self-attention mechanisms. However, due to their complexity, their latent token representations are often difficult to interpret. We introduce a novel framework that interprets Transformer embeddings, uncovering meaningful semantic patterns within them. Based on this framework, we demonstrate that zero-shot unsupervised semantic segmentation can be performed effectively without any fine-tuning using a model pre-trained for tasks other than segmentation. Our method reveals the inherent capacity of Transformer models for understanding input semantics and achieves state-of-the-art performance in semantic segmentation, outperforming traditional segmentation models. Specifically, our approach achieves an accuracy of 67.2 % and an mIoU of 32.9 % on the COCO-Stuff dataset, as well as an mIoU of 51.9 % on the PASCAL VOC dataset. Additionally, we validate our interpretability framework on LLMs for text summarization, demonstrating its broad applicability and robustness.

1. Introduction

In recent years, the Transformer architecture and foundation models, leveraging self-attention mechanisms to capture complex dependencies in text, have transformed Natural Language Processing (NLP) benchmarks [2, 51, 54, 58]. Similarly, Vision Transformers (ViTs) [13] have been adapted in Computer Vision (CV) and now serve as the backbone for various tasks such as segmentation and object detection [30, 52]. Despite their success, understanding the interpretability of Transformers remains a challenge due to the complexity of their latent token representations.

Several methods have been developed to enhance the interpretability of CNN-based models [44, 47, 63]. While some of these can be extended to Transformer architectures, they do not fully leverage the unique attention mech-

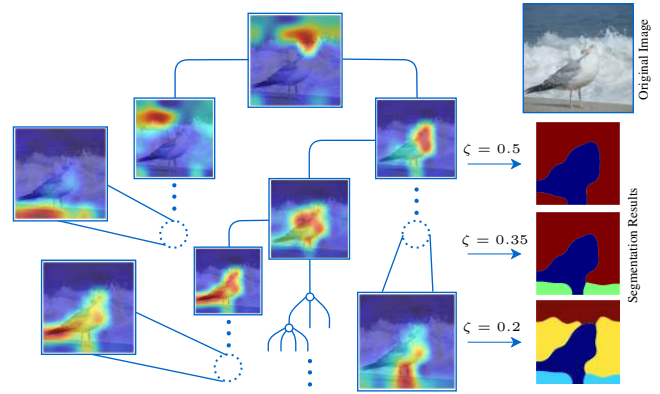


Figure 1. Hierarchical clustering tree showing the grouping of token relevance maps for all tokens in a latent layer of the Vision Transformer, not limited to the CLS token. Each leaf node represents a single token relevance map, while higher-level nodes show aggregated clusters based on a clustering threshold (ζ), which controls the level of detail. Lower ζ values reveal finer details, while higher values create broader, more general clusters. This approach demonstrates how pre-trained Vision Transformers can perform unsupervised semantic segmentation, identifying meaningful patterns within token representations without requiring additional training or fine-tuning.

anisms inherent to Transformers. Recent research has introduced interpretability methods specifically designed for Transformers [1, 6, 59]. However, these approaches primarily focus on explaining final model outputs, providing limited insight into the intermediate processes that lead to predictions. For instance, [7] maps latent tokens into CLIP’s [37] multi-modal space to find corresponding text descriptions, relying on an external text encoder for interpretability. In contrast, our approach directly interprets the latent space of ViTs, elucidating the role and function of each token within the high-dimensional space without relying on external models.

This paper introduces a framework to interpret latent to-

kens, offering a deeper understanding of the internal workings of Transformers. This understanding enables users to perform image semantic segmentation using pre-trained Transformer-based vision models in an unsupervised, zero-shot manner, without any additional training. We demonstrate that applying semantic segmentation based on our interpretability framework achieves state-of-the-art performance on benchmark image segmentation datasets.

Drawing inspiration from [6], our method analyzes the semantic information retained by latent tokens, enabling tasks such as object selection and semantic segmentation without additional training. We demonstrate that Transformers inherently understand the semantic structure of their input, viewing it as a collection of distinct concepts. Each latent token identifies a specific concept with semantic significance, thereby shedding light on the decision-making process of these models.

As shown in Section 4, our framework proves effective in a range of tasks, including semantic segmentation and model interpretation. Most recent unsupervised semantic segmentation methods involve an additional training phase to learn feature representations [18, 27, 46]. Our approach, however, utilizes the understanding embedded in pre-trained models to achieve zero-shot segmentation, leveraging their inherent knowledge of images. The stronger a model’s comprehension of image content, the more accurately it performs segmentation. We further demonstrate that our method is capable of interpreting large language models (LLMs) at the token level, validating its application in tasks such as text summarization. The main contributions of this paper are as follows:

- We propose a framework for interpreting latent tokens in Transformers, providing valuable insights into model decision-making processes.
- By aggregating relevance maps generated by tokens using hierarchical clustering, we achieve zero-shot unsupervised semantic segmentation on pre-trained models, outperforming SOTA methods that require additional training.
- We demonstrate the capability of our method to interpret LLMs at the token level, showcasing its practical applicability to textual data. Specifically, we demonstrate the interpretation of LLM operations in the text summarization task.

This paper is structured as follows: Section 2 reviews interpretability frameworks and previous work on semantic segmentation, highlighting our SOTA results. Section 3 presents our interpretability framework. In Section 4, we showcase its effectiveness on image and text tasks. Finally, Section 5 concludes the paper.

2. Related Work

The interpretability of deep learning architectures has become a central focus in AI research [16]. As models grow in complexity, understanding their decision-making processes is essential for ensuring transparency, reliability, and fairness [53]. Interpretability not only aids in debugging and performance improvement but also builds trust in AI systems, particularly in fields like healthcare, finance, and autonomous driving. Opaque models can perpetuate biases and generate unforeseen outcomes, making interpretability crucial for bridging high performance with safe, practical AI deployment [26]. In Transformer models, tokens are key to interpreting behavior, with their relationship to spatial locations in images or sequence order adding an important layer to interpretation. This relationship enhances semantic segmentation as a measure for evaluating token-based interpretation frameworks, which this section will explore.

2.1. Model Interpretation

Model interpretability is a critical area of research in deep learning, especially for complex models like transformers. Traditional deep learning models, such as CNNs, have been interpreted using various techniques like saliency maps [47], deconvolutional networks [63], Guided Backpropagation [48], Feature Visualization [34], Local Interpretable Model-Agnostic Explanations (LIME) [40], SHapley Additive exPlanations (SHAP) [31], Class Activation Mapping (CAM) [65], Feature Attribution with Integrated Gradients [50], and Grad-CAM [44], which highlight the most important regions of the input that influence model predictions. While effective for CNNs, these techniques are less suitable for transformer architectures, as they fail to account for the unique self-attention mechanisms that transformers rely on. Anchors is a model-agnostic interpretability method that provides high-precision, locally faithful explanations by generating *if-then* rules that sufficiently explain a model’s prediction for specific inputs [41].

Recent interpretability research on transformers has focused on understanding how these models allocate attention and propagate information. A seminal contribution was made by examining self-attention patterns in transformer-based language models, revealing how attention is distributed across tokens [59]. Abnar et al. [1] advanced this by proposing methods to visualize attention flow across layers, aiding in the understanding of information propagation. Chefer et al. [6] introduced a relevance-based approach, computing relevance scores for each token to provide deeper insights into the model’s decision-making process. Additionally, the sensitivity of model predictions to input tokens has been explored as another way to interpret transformer behavior [20].

Interpretability methods for reflecting semantic relations within input sequences have been investigated in NLP tasks,

aiming to characterize how transformers associate words or tokens across different positions in a sentence for downstream tasks [11]. Justifying the role of the multi-head attention mechanism in transformers and its contribution to improved performance is another area of study, revealing insights into how transformers operate [32]. Other notable interpretability methods include Interpretability-Aware Visual Transformers (IA-ViT) [36], Mechanistic Interpretability [39], Concept Transformers [42], Nested Hierarchical Transformers [64], Tracr [29], and ViT-net [22].

While many methods focus on explaining model outputs, fewer efforts have been directed toward understanding the intermediate processes within the model. The internal representations, particularly the latent tokens in Vision Transformers (ViTs), remain largely unexplored. A similar approach is the work by [7], which focuses on the CLIP model [37]. Their method involves disabling the self-attention mechanism to map the latent tokens into a multimodal space without additional training. However, simply disabling self-attention may introduce distribution shifts, raising concerns about the validity of their results. Furthermore, their approach relies heavily on CLIP’s text encoder, which limits its generalizability to other models. Our work overcomes these limitations by focusing on backward-looking analysis, investigating how latent tokens relate to the input rather than looking ahead. Additionally, we rely solely on one encoder, making our method more applicable to a broader range of models.

2.2. Unsupervised Semantic Segmentation

Unsupervised semantic segmentation has advanced significantly due to self-supervised learning and clustering-based techniques, which reduce reliance on labeled datasets. Early approaches, such as Invariant Information Clustering (IIC) [21], employed mutual information to group similar pixels without using labels. Building on this foundation, PiCIE introduced consistency by integrating photometric and geometric invariances, setting an important precedent for subsequent methods [9].

The emergence of Vision Transformers (ViTs) and self-supervised learning marked a major shift in the field. DINO became a pioneering method, extracting rich, meaningful features without the need for labeled data [5]. Its effectiveness in unsupervised segmentation laid the groundwork for more recent advancements. For instance, TransFGU [62] performs semantic segmentation in a top-down manner by deriving class activation maps from DINO models. STEGO [18] leverages DINO’s features, employing contrastive learning to group similar regions and achieving notable improvements in segmentation accuracy.

Subsequent methods further refined these concepts. MaskContrast incorporated clustering techniques to ensure region consistency across different views, enhancing feature

representations for unsupervised segmentation [56]. Leapart utilized self-supervised ViTs to improve pixel grouping, particularly in complex scenes [66]. ACSEg introduced adaptive conceptualization for pixel-level semantic grouping, using an Adaptive Concept Generator (ACG) to dynamically align learnable prototypes with relevant concepts, thereby addressing over- and under-clustering challenges in varied images [27].

Leveraging hidden positives for unsupervised semantic segmentation enhances pixel grouping by identifying and utilizing implicit positive relationships within the data, boosting segmentation performance without the need for labeled examples [45]. Using self-supervised learning and adaptive feature representations to enable models to discover and segment novel, unseen object categories without labeled data have also been investigated [55]. Smooseg introduces a smoothness prior to enhancing unsupervised semantic segmentation by promoting coherent region grouping [25]. Unsupervised semantic segmentation is also improved by leveraging depth-guided feature correlation and targeted sampling to enhance region consistency and accuracy [46]. [17].

Unlike approaches that depend on self-training, pseudo-labeling, or complex setups, our model ULTra introduces a zero-shot method for unsupervised semantic segmentation. ULTra achieves strong segmentation performance by directly leveraging the semantic information embedded in the latent tokens of pre-trained Transformer models, without the need for additional training or fine-tuning. Furthermore, ULTra emphasizes explainability within Transformers by illustrating how the model’s latent tokens contribute to segmentation, in contrast to the common practice of using only the CLS token to represent the entire image. This zero-shot, explainability-driven approach offers a novel direction for unsupervised segmentation, demonstrating that Transformer-based architectures can achieve SOTA performance without labeled data or extensive additional processing.

3. Methodology

In this section, we present our approach for interpreting latent representations in Transformers. We begin with essential preliminaries and notation for clarity, followed by a detailed explanation of our method for interpreting latent tokens and its application to semantic segmentation.

3.1. Preliminaries and Notation

The architecture of a typical Transformer can be formulated as follows: the input X is split into n tokens $\{\mathbf{x}_i\}_{i=1}^n$. After tokenization, token embeddings $\{\mathbf{e}_i\}_{i=0}^n$ are computed, where \mathbf{e}_0 corresponds to the CLS token. Positional encodings PE_i are added to the i -th token embedding to incorporate spatial information, resulting in the latent token repre-

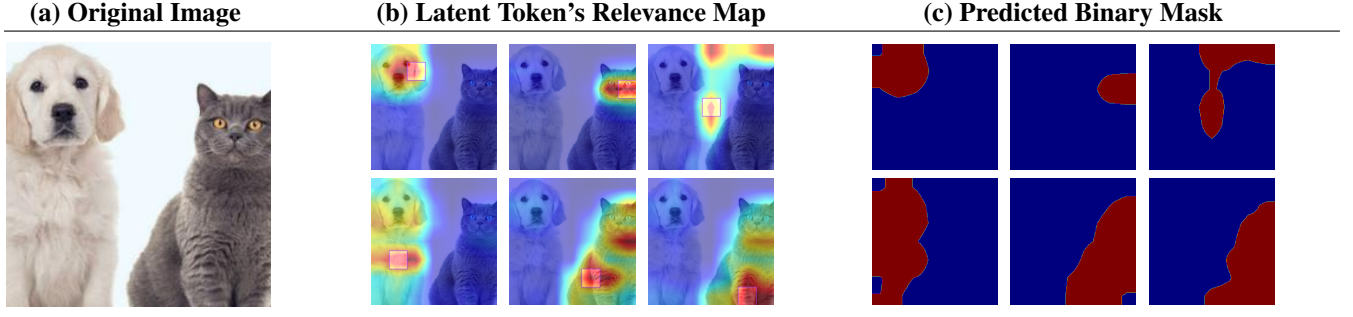


Figure 2. An example of token interpretation by our model and its predicted binary mask. (a) Original image. (b) Overlay of $\tilde{S}_i^{(13)}$ on the original image for different i , where the location of the i -th token is indicated by the purple square. (c) The binary mask $M_i^{(13)}$ for each corresponding relevance map in (b).

sentation $\mathbf{z}_i^{(1)} = \mathbf{e}_i + \text{PE}_i$. Here, $\mathbf{z}_i^{(l)}$ represents a latent token, where l denotes the layer index with $l \in \{1, \dots, L\}$ and L is the total number of layers in the Transformer, and i represents the i -th token within the l -th layer.

For each head $h \in \{1, \dots, H\}$ in the multi-head attention mechanism, the queries, keys, and values corresponding to the i -th token are obtained via linear transformations, projecting the latent token of dimension d into dimension k :

$$Q_h^{(l)}(\mathbf{z}_i^{(l-1)}) = (W_{h,q}^{(l)})^T \mathbf{z}_i^{(l-1)}, K_h^{(l)}(\mathbf{z}_i^{(l-1)}) = (W_{h,k}^{(l)})^T \mathbf{z}_i^{(l-1)},$$

$$V_h^{(l)}(\mathbf{z}_i^{(l-1)}) = (W_{h,v}^{(l)})^T \mathbf{z}_i^{(l-1)}, \quad \forall l \in \{2, \dots, L\} \quad (1)$$

where $W_{h,q}^{(l)}, W_{h,k}^{(l)}, W_{h,v}^{(l)} \in \mathbb{R}^{d \times k}$. The attention weights for each token pair (i, j) at layer l and head h are computed as:

$$\alpha_{h,i,j}^{(l)} = \text{softmax}_j \left(\frac{\langle Q_h^{(l)}(\mathbf{z}_i^{(l-1)}), K_h^{(l)}(\mathbf{z}_j^{(l-1)}) \rangle}{\sqrt{k}} \right). \quad (2)$$

Then, i -th token is updated by summing over the weighted values across all heads:

$$\bar{\mathbf{u}}_i^{(l)} = \sum_{h=1}^H (W_{c,h}^{(l)})^T \sum_{j=1}^n \alpha_{h,i,j}^{(l)} V_h^{(l)}(\mathbf{z}_j^{(l-1)}), \quad (3)$$

where $W_{c,h} \in \mathbb{R}^{k \times d}$. The updated token representation \mathbf{u}_i after the attention layer is computed as:

$$\mathbf{u}_i^{(l)} = \text{LayerNorm}(\mathbf{z}_i^{(l-1)} + \bar{\mathbf{u}}_i^{(l)}). \quad (4)$$

Each token then passes through a feed-forward network:

$$\bar{\mathbf{z}}_i^{(l)} = (W_2^{(l)})^T \text{ReLU}((W_1^{(l)})^T \mathbf{u}_i), \quad (5)$$

$$\mathbf{z}_i^{(l)} = \text{LayerNorm}(\mathbf{u}_i + \bar{\mathbf{z}}_i^{(l)}). \quad (6)$$

Here, $W_1^{(l)} \in \mathbb{R}^{d \times m}$, $W_2^{(l)} \in \mathbb{R}^{m \times d}$.

3.2. Interpreting Latent Tokens

A straightforward approach for interpreting a model involves analyzing the semantic flow from the input to the corresponding logits. This can be achieved by adding the attention probability matrix to an identity matrix I , which incorporates skip connections, and then multiplying the result across layers [1]. However, a notable challenge arises from multiple attention heads in each Transformer layer. To address this, [6] proposes performing a weighted average across attention heads, with the weights determined by the gradient of the logits with respect to the attention weights. We employed a modified, simpler version of this method, utilizing only the attention weights rather than the layer's relevance.

In our framework, latent tokens $\mathbf{z}_i^{(l)}$ are represented as vectors rather than direct class correspondences, introducing an additional layer of complexity. To manage this, an appropriate transformation, such as the energy or euclidean norm of the latent token, is employed as a surrogate. Consequently, we compute the gradient of $\|\mathbf{z}_i^{(l)}\|$ with respect to the attention weights to facilitate the analysis.

In the following equations, we define the relevance map $S_i^{(l)} \in \mathbb{R}^n$, where the j -th element of this vector represents the importance of the j -th input token to the targeted latent token $\mathbf{z}_i^{(l)}$. The relevance map is computed as:

$$\bar{\mathbf{A}}_i^{(b,l)} = I + \mathbb{E}_h \left(\nabla_{\mathbf{A}_h^b} \|\mathbf{z}_i^{(l)}\| \odot \mathbf{A}_h^{(b)} \right)^+,$$

$$\bar{S}_i^{(l)} = \bar{\mathbf{A}}_i^{(1,l)} \cdot \bar{\mathbf{A}}_i^{(2,l)} \cdot \dots \cdot \bar{\mathbf{A}}_i^{(l-1,l)},$$

$$S_i^{(l)} = \bar{S}_i^{(l)}[i, 1:], \quad (7)$$

where $(\cdot)^+$ means considering only positive values, $S, \bar{\mathbf{A}}_i^{(b,l)} \in \mathbb{R}^{(n+1) \times (n+1)}$, \odot denotes the Hadamard product, \mathbb{E}_h represents the mean across the heads dimension, and $\mathbf{A}_h^{(b)} \in \mathbb{R}^{h \times (n+1) \times (n+1)}$ is the attention score matrix in the b -th layer and $A_{h,i,j}^{(b)} = \alpha_{h,i,j}^{(b)}$.

Due to the skip connections in the transformer, most of

the contribution of $S_i^{(l)}$ is concentrated on $S_i^{(l)}[i - 1]$ which makes it hard to analyze other tokens’ contribution. To address this issue, we replace this element with the maximum value of other elements, thereby capturing the contributions of additional tokens to the selected token.

For vision tasks, we first reshape the relevance map and then upsample it using bilinear or cubic interpolation to match the resolution of the model’s input. The resulting higher-dimensional matrix is denoted as \tilde{S} . This upsampling step is essential for enabling accurate object selection and semantic segmentation tasks.

3.3. ULTra in Unsupervised Tasks

In this section, we aim to examine ULTra’s capability to adapt to various tasks involving semantic knowledge. Importantly, it requires no additional training, leveraging the inherent understanding within transformers rather than relying on loss functions objective, final layer outputs, or fine-tuning.

Unsupervised Semantic Segmentation. As previously discussed, a relevance map can be defined for each latent token at a fixed layer, with the total number of relevance maps equal to the number of latent tokens. In the context of segmentation, the goal is to assign a class label to each pixel within an image. To achieve this, we employ clustering techniques, such as hierarchical clustering, that do not require a predefined number of classes. These techniques group the relevance maps into k distinct clusters, where k is unknown. Ideally, we aim for k to approximate the actual number of classes present in the image.

Our approach provides flexibility in adjusting the value of k by modifying the cutoff distance threshold ζ within the clustering algorithm. Increasing ζ produces fewer, broader clusters that capture general categories, such as background and foreground. Conversely, reducing ζ allows for finer segmentation, distinguishing more specific features, such as an object’s head or hands. To prevent the method from disproportionately favoring larger objects, given that the number of elements in each cluster may vary, we apply min-max scaling to each cluster independently.

After clustering, we define k distinct concepts by aggregating the relevance maps within each cluster. The aggregated relevance map $\tilde{S}_c^{(l)}[x, y]$ for a cluster c is computed as:

$$\tilde{S}_{c,\zeta}^{(l)}[x, y] = \sum_{i \in \phi_\zeta(c)} \tilde{S}_i^{(l)}[x, y], \quad (8)$$

where $\phi_\zeta(c) = \{i \mid \text{Class}(\tilde{S}_i^{(l)}) = c\}$ represents the set of label assignments determined by the clustering algorithm. For each pixel at position $[x, y]$, the class label is determined by identifying the cluster c with the highest relevance value at that pixel. Mathematically, the class assignment for a pixel is expressed as:

$$\text{Class}[x, y]_\zeta^{(l)} = \underset{c \in \{1, \dots, k\}}{\text{argmax}} S_{c,\zeta}^{(l)}[x, y]. \quad (9)$$

Some examples illustrating our segmentation method is presented in Figure 3. Additionally, the hierarchy tree and the effect of the threshold are demonstrated in Figure 1.

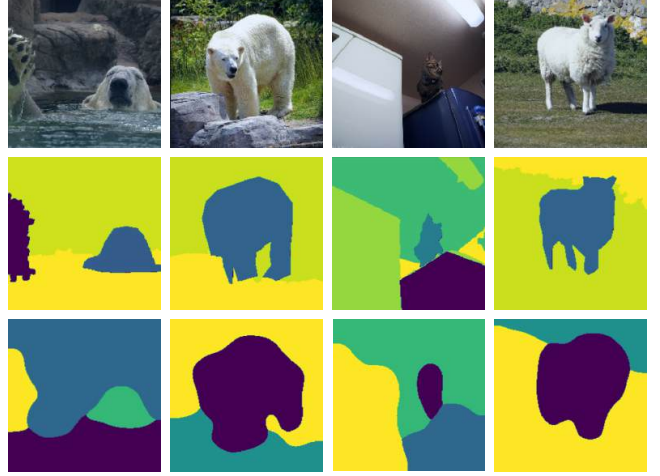


Figure 3. ULTra segmentation results on sample images. The top row displays the original images, the middle row shows ground-truth annotations, and the bottom row presents our model’s predictions.

4. Experiment

Datasets. In our experiments, we evaluate model performance on several semantic segmentation benchmarks, focusing on vision-related tasks. We conducted experiments on three datasets: COCO-Stuff 27 [4], PASCAL VOC 2012 [15], and Potsdam-3 [19]. This combination of datasets provides a diverse testing ground to rigorously evaluate our unsupervised zero-shot approach across both standard and challenging perspectives in semantic segmentation.

COCO-Stuff 27, a subset of the COCO dataset [28], includes complex, real-world scenes with pixel-level annotations across various object categories. Similarly, PASCAL VOC 2012 serves as a classic benchmark with pixel-level annotations, while the Potsdam-3 dataset offers a unique aerial, top-down perspective of urban scenes. The latter adds an additional challenge with its high-resolution images of buildings, roads, and natural elements captured over the city of Potsdam.

For our qualitative analysis of LLM interpretation in the task of text summarization, as described in Section 4.3, we utilized the TL;DR dataset [49]. The TL;DR dataset contains summary comparisons with human feedback collected by OpenAI. Each entry consists of a Reddit post, including its title, original content, and a human-generated TL;DR.

Models. For all experiments in the vision tasks, we used CLIP’s image encoder ViT-B/32 [37]. For interpreting text

summarization, as described in Section 4.3, we used the Llama-2-7B language model [54]. All experiments were run on 8 NVIDIA RTX 4090-24GB GPUs.

4.1. Zero-shot Unsupervised Object Selection

The upsampled relevance map $\tilde{S}_i^{(l)}$ can be converted into a binary segmentation mask using a threshold τ , where the binary mask $M_i^{(l)}$ is defined as:

$$M_i^{(l)}[x, y] = \begin{cases} 0, & \text{if } \tilde{S}_i^{(l)}[x, y] < \tau, \\ 1, & \text{otherwise.} \end{cases} \quad (10)$$

Here, $\tilde{S}_i^{(l)}[x, y]$ represents the relevance value at position $[x, y]$ in $\tilde{S}_i^{(l)}$, with τ as the threshold. When $\tilde{S}_i^{(l)}[x, y]$ is below τ , $M_i^{(l)}[x, y]$ is set to 0; otherwise, it is set to 1, marking the object region.

Our findings indicate that as tokens propagate through the network, they refine their object representation while retaining the semantic meaning of their associated image patches performing *Object Selection*. Figure 5 visually illustrates this process. For a given patch token \mathbf{x}_i , the object it most strongly represents is denoted as class k_i , indicating that \mathbf{x}_i predominantly corresponds to a region of the object belonging to class k_i in the image. The latent token $\mathbf{z}_i^{(l)}$ generates a relevance map that highlights areas with higher values associated with class k_i . After applying a threshold, this map becomes a binary segmentation mask expected to exhibit a high Intersection over Union (IoU) with the corresponding class k_i region in the image. An illustrative example is shown in Figure 2.

To quantify alignment, we compute the IoU by converting the relevance map $S_i^{(l)}$ into a binary mask $M_i^{(l)}$ and comparing it with the ground-truth mask. We propose the Initial Token IoU (ITIoU) metric, which measures how well the relevance maps of input tokens align with their respective class masks. The ITIoU is calculated as:

$$ITIoU^{(l)}(X) = \frac{1}{C} \sum_{i=1}^C \frac{1}{|\mathcal{P}_i|} \sum_{\mathbf{x}_j \in \mathcal{P}_i} \text{IoU}(M_j^{(l)}, G_i), \quad (11)$$

where C denotes the number of classes, \mathcal{P}_i represents the set of tokens associated with class i , $M_j^{(l)}$ is the binary segmentation mask for token \mathbf{x}_j within class i , and G_i is the ground-truth mask for class i in image X . The inner sum averages the IoU for tokens in \mathcal{P}_i for each class, and the outer sum then averages across all classes. Using a threshold of 0.2, our ITIoU metric achieves an average score of 37.84 % on the COCO validation dataset and 39.51 % on the VOC dataset.

The ITIoU metric provides a comprehensive evaluation by incorporating a weighted average across tokens in each class, enhancing the assessment of token alignment with their respective ground-truth labels.

4.2. Zero-shot Unsupervised Semantic Segmentation

We benchmarked the segmentation capability of our method against several approaches in the literature on unsupervised segmentation. Notably, unlike other methods, our approach requires no additional training. Instead, it relies solely on a pre-trained vision transformer, which may have been trained on tasks unrelated to unsupervised segmentation.

To evaluate our method, we used the Unsupervised mean Intersection over Union (U. mIoU) and Unsupervised pixel Accuracy (U. Accuracy) metrics, following standard practices in semantic segmentation research. In all experiments across datasets, as shown in Tables 1, 2, and 3, we set the cutoff distance threshold $\zeta = 0.4$.

Table 1. Comparison of unsupervised segmentation methods on the 27-class COCO-Stuff validation set. "W" indicates methods requiring additional training, while "W/O" denotes methods without additional training.

Method	Training	U. Accuracy	U. mIoU
IIC [21]	W	21.8	6.7
DINO [5]	W	30.5	9.6
PiCIE [9]	W	48.1	13.8
ACSeg [27]	W	-	16.4
STEGO [18]	W	56.9	28.2
DepthG [46]	W	58.6	29.0
U2Seg [33]	W	63.9	30.2
ULTra (Ours)	W/O	67.2	32.9

Table 2. Comparison of unsupervised segmentation methods on the PASCAL VOC validation set.

Method	Training	U. mIoU
IIC [21]	W	9.8
MaskContrast [56]	W	35.0
Leopart [66]	W	41.7
TransFGU [62]	W	37.2
MaskDistill [57]	W	42.0
ACSeg [27]	W	47.1
ULTra (Ours)	W/O	51.9

Table 3. Comparison of unsupervised segmentation methods on the Potsdam validation set.

Method	Training	U. ACC
IIC [21]	W	65.1
DINO [5]	W	71.3
STEGO [18]	W	77.0
DepthG [46]	W	80.4
ULTra (Ours)	W/O	74.6

How do I [2 0 M] stop feeling bad about myself for having no relationship experience at all ? POST : It just seems like everyone I know has at least had a " thing " with someone by this point . I ve made out with a girl once (who later told me that was a mistake) and I feel like girls always reject me or only see me as a friend . Which is perfectly acceptable , but I ' m starting to get ups et that I ' ve never had any kind of relationship . I just got rejected by a girl who I thought was into me and I ' ve been feeling bad ever since . I just don t know what ' s wrong with me . I guess I ' m a little bit skin ny (I work out regularly though) , but I show er every day , dress pretty well , all that stuff .

(a) TL;DR: I've had very bad luck with girls my whole life and I don't know how to get my confidence up.

I need help about those feelings POST : I am a 1 8 M , she ' s a 1 7 F . We ' ve got a troubles ome relationship which started as a pure friendship one year ago . I ' ve made mist akers , she made hers too . O ur last situation scenario is explained in here : Now I feel like I hate her , I used to adm ire her a lot , but I ' m really disappoint ed with her and with her character . But I just realized I still like her . So , well , yeah , I like her and hate her . And just after that bad situation happened I realized she also had that feeling . Well , now we both hate and love each other . What to do ? What to think ? What to feel ? add itional info : today our friend asked me for help with some calculations and I made a jo ke about our physics teacher . She laughed and smiled at me just like one year ago , but after she realized that , she seemed kind a [gr ouch y] .

(b) TL;DR: I still like her but my rational side says "no, she is a trash person".

Figure 4. Visualization of Token Contribution Scores ($\lambda_i^{(l)}$) highlighting the relevance of context tokens in interpreting the summary. Each token is colored proportionally to its $\lambda_i^{(l)}$ value. These visualizations demonstrate the model's ability to identify key semantic elements in the context for generating relevant summaries.

Our results in Tables 1, 2, and 3 highlight the effectiveness of the proposed method for unsupervised semantic segmentation. Unlike baseline approaches, which require additional training specifically tailored to segmentation, our method leverages a pre-trained vision transformer without any further fine-tuning. Despite this lack of task-specific training, our approach consistently outperforms other methods, demonstrating adaptability across diverse datasets.

4.3. Interpreting LLMs in Text Summarization

In this section, we examine how our interpretability framework can be applied to text summarization tasks to uncover the underlying mechanisms and intent of LLMs. By visualizing the regions of the input context that an LLM prioritizes while interpreting a given TL;DR summary, we gain deeper insights into the model's behavior and decision-making processes. As shown in Figure 4, these visualizations allow us to pinpoint the most influential regions of a textual input prompt in generating concise and relevant summaries.

For the experiments, we used a Supervised Fine-Tuned (SFT) version of Llama-2-7B trained on the UltraFeedback Binarized (UFB) dataset [12]. Additionally, we aligned the model to the text summarization task on the TL;DR dataset using the Direct Preference Optimization (DPO) method [38] for 1,000 iterations, with a learning rate of 5×10^{-6} and $\beta = 0.5$. To validate our framework, we selected the preferred response (TL;DR) of each sample in the dataset, denoted it as y , and used it as the summary of the context x .

In this task, we concatenate the context x and the summary y with a separator token. After feeding this input to the model, we compute the relevance scores of the TL;DR tokens with respect to the context tokens. We then average these scores for each token in x to obtain a scalar value, referred to as the Token Contribution Score, $\lambda_i^{(l)} \in \mathbb{R}^+$, which highlights the contribution of each context token in interpreting the summary y with respect to the context. This provides visual evidence of how the model processes the context text and identifies key semantic elements relevant

to producing the summary y . Accordingly, $\lambda_i^{(l)}$ is computed as:

$$\lambda_i^{(l)} = \frac{1}{|y|} \sum_{j=1}^{|y|} S_{j+|x|}^{(l)}[i], \quad \forall i \in \{1, \dots, |x|\}, \quad (12)$$

where $|\cdot|$ denotes the number of tokens in the text.

In example (a): semantically significant words such as 'relationship', 'experience', 'rejection', and 'never' are prominently highlighted, reflecting the model's interpretation of the person's struggles with relationships and feelings of rejection. Additionally, the highlighting of the question at the beginning of the context 'How do I stop feeling bad...' suggests the model recognizes the presence of uncertainty and a request for guidance, which is encapsulated in the summary as 'I don't know.'

In example (b): $\lambda_i^{(l)}$ scores reveals the model's focus on words such as 'feelings', 'hate', 'disappoint', 'love', and 'like', which correspond to the person's mixed emotions toward their girlfriend, as described in the summary. The apparent contradiction between 'love' and 'trashness' in the summary seems to be derived from these highlighted terms, suggesting the model understands the conflicting emotions present in the text. Furthermore, the emphasis on 'character' aligns with the summary's judgmental tone, suggesting that the model interprets this word as indicative of an assessment of personality traits and behaviors.

This token-level analysis can serve as a valuable tool for effective supervision in future research, particularly for developing interpretable fine-tuning and alignment methods like Reinforcement Learning from Human Feedback (RLHF) [10, 35, 49] and Preference Optimization [3, 14, 38], where understanding model behavior and intent is essential. It may also provide useful insights in related areas, such as Chain-of-Thought [8, 24, 43] and Theory-of-Mind [23, 60, 61], which seek to make the reasoning processes and intentions of LLMs more transparent. Ultimately, our framework offers insights into fundamental questions in

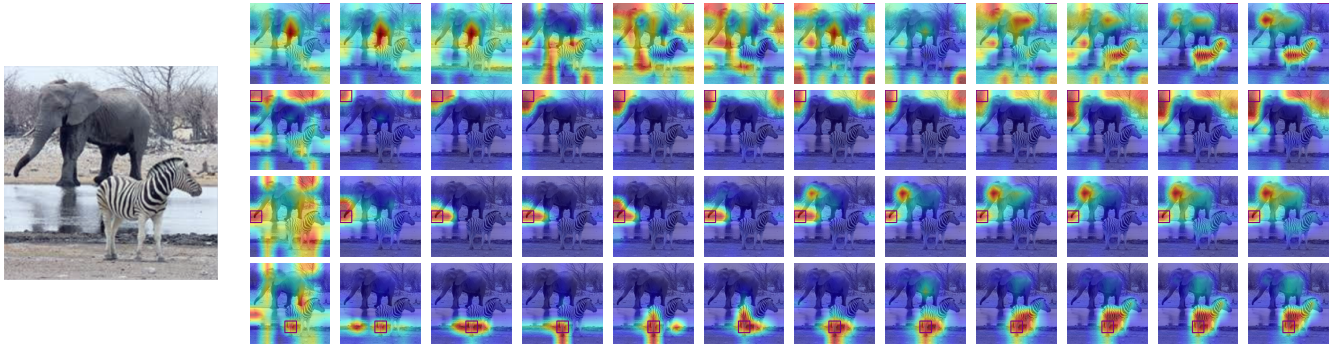
(a) Original Image**(b) Explanation Relevancy Map**

Figure 5. An example illustrating the model’s decision-making process across layers. Moving from left to right corresponds to deeper layers in the network. The first row corresponds to the CLS token, while the second, third, and fourth rows represent three different tokens, highlighted by red squares.

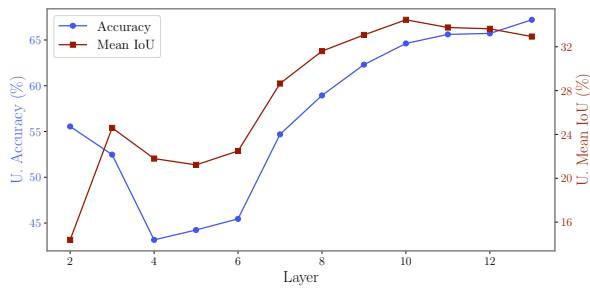
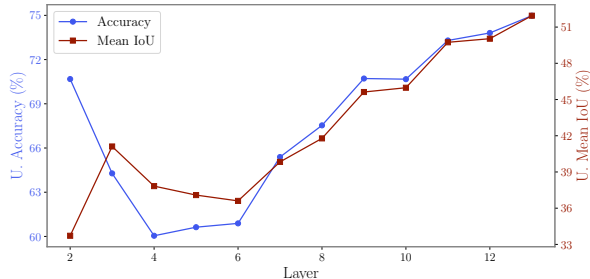
**(a) COCO-Stuff****(b) PASCAL VOC**

Figure 6. Ablation study on two evaluation metrics across layers. These plots demonstrate a progressive improvement in semantic segmentation performance in the deeper layers of the transformer model. This enhancement is attributed to latent tokens capturing more meaningful segment structures, resulting in increasingly accurate and refined semantic representations.

NLP, such as: *“How is the model thinking?”* or *“What is the underlying intent behind the model’s generated response?”*.

4.4. Ablation Study on the Effect of Layer Depth in ViT Token Understanding

In this section, we analyze the impact of depth on our model’s interpretability and segmentation performance,

providing insights into the contribution of each layer. As anticipated and observed in Figure 6, deeper layers generally carry more semantic significance. However, the contribution diminishes in the final layers, suggesting that a depth of around 13 layers might be more than sufficient for the ViT to effectively comprehend image content. This finding implies that even fewer layers might achieve comparable results, potentially reducing computational costs without compromising performance.

We observe an intriguing behavior in the initial layers, where performance initially declines before improving. This phenomenon is also visually evident in Figure 5, where the attention maps in the first layer appear to focus on the entire image. This suggests that, initially, the token examines the image as a whole before selectively gathering information from tokens with similar characteristics.

5. Conclusion

This paper introduces a novel interpretability framework that provides valuable insights into the decision-making processes of Transformer models. The framework is based on the input interpretation provided by each token in an arbitrary layer. By aggregating the interpretation of all tokens in the same layer, the framework enhances understanding of the Transformer’s behavior at the layer level. The richness of this understanding enables zero-shot unsupervised semantic segmentation, where our method achieves state-of-the-art performance with an accuracy of 67.2% and an mIoU of 32.9% on the COCO-Stuff dataset, and an mIoU of 51.9% on the PASCAL VOC dataset for the unsupervised zero-shot version. Beyond vision tasks, the framework is not limited to vision models and can also be applied to interpret the behavior of large language models (LLMs) in tasks such as text summarization. This highlights the versatility and broad applicability of the proposed framework

across various domains, offering potential for future work to simultaneously demystify the multimodal understanding of Transformers.

References

- [1] Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. *arXiv preprint arXiv:2005.00928*, 2020. 1, 2, 4
- [2] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L el io Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. Mistral 7B, 2023. arXiv:2310.06825 [cs]. 1
- [3] Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and R emi Munos. A general theoretical paradigm to understand learning from human preferences, 2023. 7
- [4] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocomstuff: Thing and stuff classes in context. In *Computer vision and pattern recognition (CVPR), 2018 IEEE conference on*. IEEE, 2018. 5
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Herv e J egou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3, 6
- [6] Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 782–791, 2021. 1, 2, 4
- [7] Haozhe Chen, Junfeng Yang, Carl Vondrick, and Chengzhi Mao. Invite: Interpret and control vision-language models with text explanations. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 3
- [8] Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yungwei Lai, Zexuan Xiong, and Minlie Huang. Tombench: Benchmarking theory of mind in large language models, 2024. 7
- [9] Jang Hyun Cho, Utkarsh Mall, Kavita Bala, and Bharath Hariharan. Picie: Unsupervised semantic segmentation using invariance and equivariance in clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16794–16804, 2021. 3, 6
- [10] Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences, 2023. 7
- [11] Kevin Clark. What does bert look at? an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341*, 2019. 3
- [12] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with scaled ai feedback, 2024. 7
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1
- [14] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization, 2024. 7
- [15] Mark Everingham and John Winn. The pascal visual object classes challenge 2011 (voc2011) development kit. *Pattern Analysis, Statistical Modelling and Computational Learning, Tech. Rep.*, 8, 2011. 5
- [16] Lei Gao and Ling Guan. Interpretability of machine learning: Recent advances and future prospects. *IEEE MultiMedia*, 30(4):105–118, 2023. 2
- [17] Shanghua Gao, Zhong-Yu Li, Ming-Hsuan Yang, Ming-Ming Cheng, Junwei Han, and Philip Torr. Large-scale unsupervised semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 45(6):7457–7476, 2022. 3
- [18] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snaveley, and William T Freeman. Unsupervised semantic segmentation by distilling feature correspondences. *arXiv preprint arXiv:2203.08414*, 2022. 2, 3, 6
- [19] ISPRS. 2d semantic labeling contest - potsdam. <https://www.isprs.org/education/benchmarks/UrbanSemLab/2d-sem-label-potsdam.aspx>, 2018. 5
- [20] Sarthak Jain and Byron C Wallace. Attention is not explanation. *arXiv preprint arXiv:1902.10186*, 2019. 2
- [21] Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9865–9874, 2019. 3, 6
- [22] Sangwon Kim, Jaeyeal Nam, and Byoung Chul Ko. Vit-net: Interpretable vision transformers with neural tree decoder. In *International conference on machine learning*, pages 11162–11172. PMLR, 2022. 3
- [23] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023. 7
- [24] Michal Kosinski. Evaluating large language models in theory of mind tasks. *Proceedings of the National Academy of Sciences*, 121(45), 2024. 7
- [25] Mengcheng Lan, Xinjiang Wang, Yiping Ke, Jiaying Xu, Litong Feng, and Wayne Zhang. Smooseg: smoothness prior for unsupervised semantic segmentation. *Advances in Neural Information Processing Systems*, 36:11353–11373, 2023. 3
- [26] Bo Li, Peng Qi, Bo Liu, Shuai Di, Jingen Liu, Jiquan Pei, Jinfeng Yi, and Bowen Zhou. Trustworthy ai: From principles to practices. *ACM Computing Surveys*, 55(9):1–46, 2023. 2
- [27] Kehan Li, Zhennan Wang, Zesen Cheng, Runyi Yu, Yian Zhao, Guoli Song, Chang Liu, Li Yuan, and Jie Chen. Acseg: Adaptive conceptualization for unsupervised semantic segmentation. In *Proceedings of the IEEE/CVF conference on*

- computer vision and pattern recognition*, pages 7162–7172, 2023. 2, 3, 6
- [28] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5
- [29] David Lindner, János Kramár, Sebastian Farquhar, Matthew Rahtz, Tom McGrath, and Vladimir Mikulik. Tracr: Compiled transformers as a laboratory for interpretability. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [30] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 1
- [31] Scott Lundberg. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*, 2017. 2
- [32] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019. 3
- [33] Dantong Niu, Xudong Wang, Xinyang Han, Long Lian, Roei Herzig, and Trevor Darrell. Unsupervised universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22744–22754, 2024. 6
- [34] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11):e7, 2017. 2
- [35] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. 7
- [36] Yao Qiang, Chengyin Li, Prashant Khanduri, and Dongxiao Zhu. Interpretability-aware vision transformer. *arXiv preprint arXiv:2309.08035*, 2023. 3
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 3, 5
- [38] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. 7
- [39] Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646*, 2024. 3
- [40] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ” why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016. 2
- [41] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 2
- [42] Mattia Rigotti, Christoph Mikšovic, Ioana Giurgiu, Thomas Gschwind, and Paolo Scotton. Attention-based interpretability with concept transformers. In *International conference on learning representations*, 2021. 3
- [43] Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. Minding language models’ (lack of) theory of mind: A plug-and-play multi-character belief tracker, 2023. 7
- [44] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 1, 2
- [45] Hyun Seok Seong, WonJun Moon, SuBeen Lee, and Jae-Pil Heo. Leveraging hidden positives for unsupervised semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19540–19549, 2023. 3
- [46] Leon Sick, Dominik Engel, Pedro Hermosilla, and Timo Ropinski. Unsupervised semantic segmentation through depth-guided feature correlation and sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3637–3646, 2024. 2, 3, 6
- [47] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2014. 1, 2
- [48] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 2
- [49] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022. 5, 7
- [50] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017. 2
- [51] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste

- Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Ser-tan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. 1
- [52] Hans Thisanke, Chamli Deshan, Kavindu Chamith, Sachith Seneviratne, Rajith Vidanaarachchi, and Damayanthi Herath. Semantic segmentation using vision transformers: A survey. *Engineering Applications of Artificial Intelligence*, 126:106669, 2023. 1
- [53] Bhavani Thuraisingham. Trustworthy machine learning. *IEEE Intelligent Systems*, 37(1):21–24, 2022. 2
- [54] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutí Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. 1, 6
- [55] Svenja Uhlemeyer, Matthias Rottmann, and Hanno Gottschalk. Towards unsupervised open world semantic segmentation. In *Uncertainty in Artificial Intelligence*, pages 1981–1991. PMLR, 2022. 3
- [56] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Unsupervised semantic segmentation by contrasting object mask proposals. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10052–10062, 2021. 3, 6
- [57] Wouter Van Gansbeke, Simon Vandenhende, and Luc Van Gool. Discovering object masks with transformers for unsupervised semantic segmentation. *arXiv preprint arXiv:2206.06363*, 2022. 6
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 1
- [59] Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model. *arXiv preprint arXiv:1906.04284*, 2019. 1, 2
- [60] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. 7
- [61] Yu Xia, Rui Wang, Xu Liu, Mingyan Li, Tong Yu, Xiang Chen, Julian McAuley, and Shuai Li. Beyond chain-of-thought: A survey of chain-of-x paradigms for llms, 2024. 7
- [62] Zhaoyuan Yin, Pichao Wang, Fan Wang, Xianzhe Xu, Hanling Zhang, Hao Li, and Rong Jin. Transfgu: a top-down approach to fine-grained unsupervised semantic segmentation. In *European conference on computer vision*, pages 73–89. Springer, 2022. 3, 6
- [63] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. *arXiv preprint arXiv:1311.2901*, 2014. 1, 2
- [64] Zizhao Zhang, Han Zhang, Long Zhao, Ting Chen, Ser-can Ö Arik, and Tomas Pfister. Nested hierarchical transformer: Towards accurate, data-efficient and interpretable visual understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3417–3425, 2022. 3
- [65] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. 2
- [66] Adrian Ziegler and Yuki M Asano. Self-supervised learning of object parts for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14502–14511, 2022. 3, 6