

CLIP-Event: Connecting Text and Images with Event Structures

Manling Li^{1*}, Ruochen Xu², Shuohang Wang², Luowei Zhou², Xudong Lin³
Chenguang Zhu², Michael Zeng², Heng Ji¹, Shih-Fu Chang³

¹University of Illinois Urbana-Champaign ²Microsoft Research ³Columbia University
{manling2, hengji}@illinois.edu,
{ruox, shuowa, luozhou, chezhu, nzeng}@microsoft.com
{xudong.lin, sc250}@columbia.edu

Abstract

Vision-language (V+L) pretraining models have achieved great success in supporting multimedia applications by understanding the alignments between images and text. While existing vision-language pretraining models primarily focus on understanding objects in images or entities in text, they often ignore the alignment at the level of events and their argument structures. In this work, we propose a contrastive learning framework to enforce vision-language pretraining models to comprehend events and associated argument (participant) roles. To achieve this, we take advantage of text information extraction technologies to obtain event structural knowledge, and utilize multiple prompt functions to contrast difficult negative descriptions by manipulating event structures. We also design an event graph alignment loss based on optimal transport to capture event argument structures. In addition, we collect a large event-rich dataset (106,875 images) for pretraining, which provides a more challenging image retrieval benchmark to assess the understanding of complicated lengthy sentences¹. Experiments show that our zero-shot CLIP-Event outperforms the state-of-the-art supervised model in argument extraction on Multimedia Event Extraction, achieving more than 5% absolute F-score gain in event extraction, as well as significant improvements on a variety of downstream tasks under zero-shot settings.

1. Introduction

Real-world multimedia applications require an understanding of not only entity knowledge (i.e., objects and object types), but also event knowledge (i.e., event types) with event argument structures (i.e., entities involved and their roles). For example, 89% images include events in contem-



Figure 1. Examples of visual event ATTACK with different arguments. Groundings are bounding-boxes colored to match roles.

porary multimedia news data². Furthermore, recognizing the arguments (participants) is critical for news comprehension, since events might be contradictory if the arguments play different roles. For example, both Fig. 1(a) and Fig. 1(b) are about the same event type ATTACK and contain entities *protester* and *police*, but their argument roles are different, i.e., the *protester* plays the role of ATTACKER in the first event and the role of TARGET in the second event, and vice versa for the *police*. Different argument roles for the same group entity result in the differentiation of two attack events.

However, existing vision-language pretraining models [4, 12, 18, 26, 31, 41] focus on the understanding of images or entities, ignoring the event semantics and structures. As a result, apparent failures happen in the circumstances requiring verb comprehension [10]. Thus, we focus on integrating event structural knowledge into vision-language pretraining. Previous work primarily represents visual events as verbs with subjects and objects [13, 19, 30, 33, 36, 43]. However, events contain structural knowledge, with each event being assigned to an *event type* that represents a set of synonymous verbs. Each argument is grounded to text or images, and associated with an *argument role* that the participant is playing. As shown in Fig. 2, the *carry* event is typed as TRANSPORT, with *protesters* as AGENT, *injured man* as

* The work is done when the first author was an intern at Microsoft.

¹ The data and code are publicly available for research purpose in <https://github.com/limanling/clip-event>.

²We randomly check 100 images at <https://www.voanews.com/>.

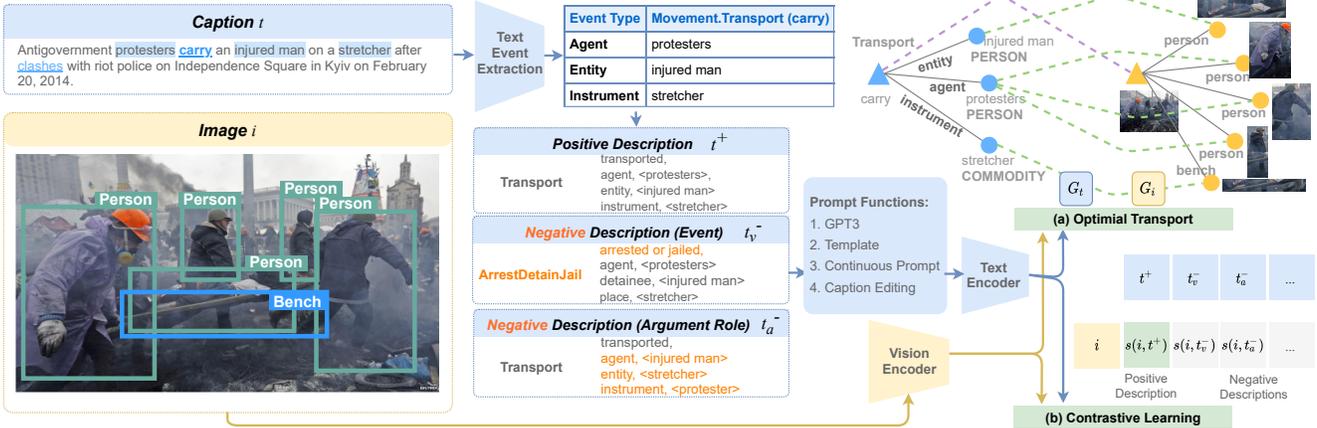


Figure 2. Architecture of CLIP-Event. We take advantage of event structural knowledge in captions to contrast hard negatives about event types and argument roles (in blue), which is then used to supervise image event understanding (in yellow) as a cross-media transfer of event knowledge. The negative event structures are highlighted in orange. The events and objects are from automatic system output.

ENTITY and *stretcher* as INSTRUMENT.

There has been little research [17, 25] on extracting event structures from news images, yielding limited support for event knowledge acquisition needed in downstream applications. Thus, we propose to leverage text information extraction technologies, which have been well researched in natural language processing, to automatically extract event structures from captions. The captions essentially refer to the same event as the images in news data, e.g., 87% captions describe the events in the images³. Therefore, we design a self-supervised contrastive learning framework, **CLIP-Event**, using the rich event knowledge in captions as distant supervision to interpret events in the associated images, to effectively transfer event knowledge across modalities.

In addition, in order to train robust representations capable of discriminating subtle differences between event types (e.g. TRANSPORT and ARREST) and argument roles (e.g. ATTACKER and VICTIM) using only images, we propose to generate *hard negatives* by manipulating event structures. We translate both correct and manipulated event structures into text descriptions using an extensive set of *event prompt functions*. Following the state-of-the-art vision-language pre-training model CLIP [26], we optimize a contrastive learning objective between images and event-aware text descriptions.

Furthermore, to transfer knowledge of argument structures, we explicitly construct event graphs consisting of event types and argument roles in vision and text. We introduce a *fine-grained* alignment between two event graphs, aligning the objects in images with the corresponding text entities and their argument roles. We employ optimal transport to encourage a *global* alignment based on the structures of two graphs, which enables the model to capture the interactions

between arguments. For example, objects with similar visual features tend to be aligned to the same argument role.

Our evaluations mainly focus on zero-shot settings, since it is crucial to understand new or previously unidentified events in real-world applications. Traditional methods based on limited pre-defined event ontologies are inapplicable in dealing with open-world events. Our pretrained model, on the other hand, is able to identify event structure using the natural language description of any unseen type and argument role, enabling zero-shot multimedia event extraction.

The evaluations on Multimedia Event Extraction [17] and Grounded Situation Recognition [25] show that CLIP-Event significantly outperforms state-of-the-art vision-language pretraining models under both zero-shot settings and supervised settings. Furthermore, it achieves significant gains in various downstream tasks under zero-shot settings such as image retrieval [7], visual commonsense reasoning [40] and visual commonsense reasoning in time [24].

In summary, this paper makes the following contributions:

- We are the first to exploit the visual event and argument structure information in vision-language pretraining.
- We introduce a novel framework by contrasting with negative event descriptions, which are generated by various prompt functions conditioned on hard negative events and arguments.
- We propose event graph alignment based on optimal transport, extending previous image or object alignment to event structure aware alignment.
- We release an event-rich image-caption dataset with 106,875 images, including the extracted event knowledge, which can serve as a challenging image retrieval

³The statistics are on those mentioned above 100 randomly sampled images from VOA News [1].

benchmark for evaluating the ability to understand complex and lengthy sentences in real-world applications.

2. Our Approach

Our goal is to incorporate event structured knowledge into vision-language pretraining. In the following we will address two primary questions regarding the model design: (1) How can the structural event knowledge be acquired? (2) How can the semantics and structures of events be encoded? We define the symbols used in this paper in Tab. 2.

2.1. Event Structural Knowledge Extraction

Text and Visual Knowledge Extraction. We use a state-of-the-art text information extraction system [16, 20] to extract events of 187 types⁴, covering a wide range of newsworthy events. For images, we apply Faster R-CNN [27] trained on Open Images [15] to detect objects.

Primary Event Detection. When there are multiple events in the caption, the image typically depicts the primary event of the caption. We detect the primary event as the event that is closer to the root of dependency parsing tree [23], and has a larger number of arguments, higher event type frequency, and higher similarity between trigger word and the image using the pretrained CLIP model [26]. We rank events according to these criteria, and perform majority voting. For example, in Fig. 2, there are two events *carry* and *clashes* in the caption. We select *carry* as the primary event since it is the root of the dependency tree, and it has three arguments, as well as higher similarity with the image.

2.2. Event Structure Driven Negative Sampling

To force the Text and Vision Encoders to learn robust features about event types and argument roles, we design the following strategies to generate challenging negatives.

Negative Event Sampling. We compute the confusion matrix for the event type classifier of the state-of-the-art vision-language pretraining model CLIP [26] on the pretraining image-caption dataset. The classifier is based on the similarity scores between the event type labels $\phi_v \in \Phi_V$ (such as TRANSPORT) and the input image i , and select the top one as the predicted event type ϕ_v^* .

$$\phi_v^* = \arg \max_{\phi_v \in \Phi_V} \phi_v^T \cdot i,$$

where the bold symbols stand for the representations from the Text and Vision Encoders in Fig. 2, and we follow CLIP to use Text and Vision Transformers. The confusion matrix is computed by comparing the predicted event type with the type of the primary event for the image. As a result, the negative event types are the challenging cases in image event

⁴The system uses DARPA AIDA ontology, which is the most fine-grained text event ontology, as attached in the Appendix.

typing, i.e., the event types whose visual features are ambiguous with the primary event type. For example, in Fig. 2, ARREST is sampled as a negative event type, since its visual features are similar to TRANSPORT.

Negative Argument Sampling. For argument roles, since each event by definition has multiple arguments, we manipulate the order of arguments by performing a right-rotation of the argument role sequence. In detail, we first order existing argument roles following the ontology definition, such as “AGENT, ENTITY, INSTRUMENT” in Fig. 2. After that, we right rotate the argument role sequence by one step, resulting in “INSTRUMENT, AGENT, ENTITY”. As a result, each argument is re-assigned to a manipulated role, e.g., *injured man*, the second argument, is manipulated from ENTITY to AGENT. If there is only one argument for the event, we sample a negative role according to the argument confusion matrix of the text argument extraction system [20].

Description Generation. To encode the positive and negative event structures using the Text Encoder, we design multiple prompt functions, as shown in Tab. 1: (1) **Single Template-based Prompt** encodes all arguments in one sentence. (2) **Composed Template-based Prompt** uses a short sentence to each argument. (3) **Continuous Prompt** employs learnable prepended tokens $[X_i]$. (4) **Caption Editing** has minimum information loss by only altering event trigger word or switching arguments. (5) **GPT-3 based Prompt** generates a semantically coherent natural language description conditioned on the event structure. We employ GPT-3 [8] and use five manual event description examples as few-shot prompts [8] to control the generation. The input to GPT-3 is the concatenation of the example events ($[ex_v]$) with arguments ($[ex_a]$), the example descriptions ($[ex_desp]$), and the target events ($[input_v]$) with arguments ($[input_a]$). The output of GPT-3 is the target description ($[output_desp]$). The description is more natural compared to template-based methods.

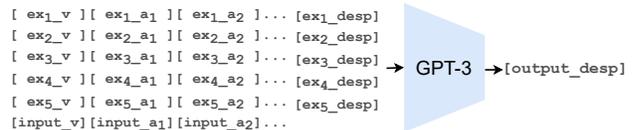


Figure 3. Architecture of GPT-3 based prompt.

2.3. Event Graph Alignment via Optimal Transport

Each event and its arguments can be organized as a graph, as shown in Fig. 2, where the central node is the event node (triangle nodes), and it’s connected to entities (circle nodes) via argument roles. Encoding event graph structures enables the model to capture the interactions between events and arguments. For example, the *injured man* should be aligned with the ENTITY being transported, rather than the AGENT.

Prompt	Example descriptions of Fig. 2 with arrest as negative event	
Single Template	Template	$\langle \text{arg1} \rangle$ transported $\langle \text{arg2} \rangle$ in $\langle \text{arg3} \rangle$ instrument from $\langle \text{arg4} \rangle$ place to $\langle \text{arg5} \rangle$ place.
	Positive	<i>Protesters</i> transported <u><i>an injured man</i></u> in <u><i>a stretcher</i></u> instrument.
	Negative-Evt	<i>Protesters</i> arrested <u><i>an injured man</i></u> in <u><i>a stretcher</i></u> place.
	Negative-Arg	<u><i>An injured man</i></u> transported <u><i>a stretcher</i></u> in <u><i>protesters</i></u> instrument.
Composed Template	Template	The image is about TRANSPORT . The AGENT is $\langle \text{arg1} \rangle$. The ENTITY is $\langle \text{arg2} \rangle$. The INSTRUMENT in $\langle \text{arg3} \rangle$. The ORIGIN is $\langle \text{arg4} \rangle$. The DESTINATION is $\langle \text{arg5} \rangle$.
	Positive	The image is about TRANSPORT . The AGENT is <u><i>protesters</i></u> . The ENTITY is <u><i>an injured man</i></u> . The INSTRUMENT is <u><i>a stretcher</i></u> .
	Negative-Evt	The image is about ARREST . The AGENT is <u><i>protesters</i></u> . The DETAINEE is <u><i>an injured man</i></u> . The PLACE is <u><i>a stretcher</i></u> .
	Negative-Arg	The image is about TRANSPORT . The AGENT is <u><i>an injured man</i></u> . The ENTITY is <u><i>a stretcher</i></u> . The INSTRUMENT is <u><i>protesters</i></u> .
Continuous Prompt	Template	$[X_0]$ TRANSPORT $[X_1]$ AGENT $[X_2]$ $\langle \text{arg1} \rangle$ $[X_3]$ ENTITY $[X_2]$ $\langle \text{arg2} \rangle$ $[X_3]$ INSTRUMENT $[X_2]$ $\langle \text{arg3} \rangle$ $[X_3]$ ORIGIN $[X_2]$ $\langle \text{arg4} \rangle$ $[X_3]$ DESTINATION $[X_2]$ $\langle \text{arg5} \rangle$ $[X_3]$
	Positive	$[X_0]$ TRANSPORT $[X_1]$ AGENT $[X_2]$ <u><i>protesters</i></u> $[X_3]$ ENTITY $[X_2]$ <u><i>an injured man</i></u> $[X_3]$ INSTRUMENT $[X_2]$ <u><i>a stretcher</i></u> $[X_3]$
	Negative-Evt	$[X_0]$ ARREST $[X_1]$ AGENT $[X_2]$ <u><i>protesters</i></u> $[X_3]$ DETAINEE $[X_2]$ <u><i>an injured man</i></u> $[X_3]$ PLACE $[X_2]$ <u><i>a stretcher</i></u> $[X_3]$
	Negative-Arg	$[X_0]$ TRANSPORT $[X_1]$ AGENT $[X_2]$ <u><i>an injured man</i></u> $[X_3]$ ENTITY $[X_2]$ <u><i>a stretcher</i></u> $[X_3]$ INSTRUMENT $[X_2]$ <u><i>protesters</i></u> $[X_3]$
Caption Editing	Positive	<i>Antigovernment protesters</i> carry <u><i>an injured man</i></u> on <u><i>a stretcher</i></u> after clashes with riot police on Independence Square in ...
	Negative-Evt	<i>Antigovernment protesters</i> arrest <u><i>an injured man</i></u> on <u><i>a stretcher</i></u> after clashes with riot police on Independence Square in ...
	Negative-Arg	<u><i>An injured man</i></u> carry <u><i>a stretcher</i></u> on <u><i>antigovernment protesters</i></u> after clashes with riot police on Independence Square in ...
GPT-3	Positive	<i>Protesters</i> transported <u><i>an injured man</i></u> with <u><i>a stretcher</i></u> .
	Negative-Evt	<i>Protesters</i> arrested <u><i>an injured man</i></u> with <u><i>a stretcher</i></u> .
	Negative-Arg	<u><i>An injured man</i></u> transported <u><i>a stretcher</i></u> and <u><i>protesters</i></u> .

Table 1. The automatically generated positive and negative descriptions for Fig. 2. We use **bold** to represent events, and underline stands for arguments. The corrupted event type and arguments are in **orange**, and templates are in **blue**. $[X_i]$ is learnable prepended token.

Symbol	Meaning
$\langle i, t \rangle$	image i and its caption text t
o, ϕ_o, i_o	object, object type, object bounding box
e, ϕ_e, t_e	entity, entity type, entity text mention
v, ϕ_v, t_v	event, event type, event text mention
$a \in \mathcal{A}(v)$	argument role; $\mathcal{A}(v)$ is the Argument role set of event v , defined by the IE ontology ³
G_i, G_t	event graph from image i and text t
t^+, t_v^-, t_a^-	positive description, negative event description, negative argument description

Table 2. List of symbols.

1. Image-level Alignment. We compute cosine similarity $s(t, i)$ and distance $d(t, i)$ between the text t and image i :

$$s(t, i) = \cos(\mathbf{t}, \mathbf{i}), d(t, i) = c(\mathbf{t}, \mathbf{i}),$$

where $c(\cdot, \cdot) = 1 - \cos(\cdot, \cdot)$ is the cosine distance function, and \mathbf{t} is obtained from the Text Transformer and \mathbf{i} is obtained from the Vision Transformer.

2. Entity-level Alignment. The cosine distance between text entity e and image object o considers both the mention

similarity and type similarity.

$$d(e, o) = c(\mathbf{t}_e, \mathbf{i}_o) + c(\phi_e, \phi_o),$$

where t_e is the text mention of entity e , and \mathbf{t}_e is its embedding contextualized on the sentence. We encode the sentence using the Text Transformer following [26], and apply average pooling over the tokens in the entity mention t_e . Similarly, i_o is the bounding box of object o and \mathbf{i}_o is its embedding contextualized on the image, based on the average pooling over the Vision Transformer representations of the patches covered in the bounding box. ϕ_e and ϕ_o are the type representations encoded by the Text Transformer. For example, $\phi_e = \text{PERSON}$ for $e = \textit{injured man}$ and $\phi_o = \text{PERSON}$ for $o = \text{[img]}$. Therefore, the distance between the aforementioned entity and object is:

$$d(e, o) = c(\textit{injured man}, \text{[img]}) + c(\text{PERSON}, \text{PERSON}),$$

3. Event-level Alignment. To obtain a global alignment score based on the structures of two graphs, we use the optimal transport [29] to get the minimal distance $d(G_t, G_i)$ between text event graph G_t and image event graph G_i ,

$$d(G_t, G_i) = \min_{\mathbf{T}} \mathbf{T} \odot \mathbf{C},$$

where \odot represents the Hadamard product. $\mathbf{T} \in \mathbb{R}_+^{n \times m}$ denotes the transport plan, learned to optimize a *soft* node alignment between two graphs. n and m are the numbers of nodes in G_t and G_i , respectively. Namely, each node in text graph G_t can be transferred to multiple nodes in image graph G_i with different weights.

C is the cost matrix. We define cost between event nodes, and between argument nodes. For event nodes, the cost is the cosine distance between the image i and trigger word v ,

$$C(v, i) = c(\mathbf{t}_v, \mathbf{i}) + c(\phi_v, \mathbf{i}).$$

For example, in Fig. 2, $v = \text{carry}$ and $\phi_v = \text{TRANSPORT}$,

$$C(v, i) = c(\text{carry}, \text{img}) + c(\text{TRANSPORT}, \text{img}).$$

The representation \mathbf{t}_v is also from the Text Transformer, contextualized on the text sentence.

The cost between each argument $\langle a, e \rangle$ and each bounding box o is based on the similarity of object o with both argument role a and text entity e .

$$\begin{aligned} C(\langle a, e \rangle, o) &= d(a, o) + d(e, o) \\ &= c(\mathbf{t}_a, \mathbf{i}_o) + c(\mathbf{t}_e, \mathbf{i}_o) + c(\phi_e, \phi_o), \end{aligned}$$

where t_a is the argument description. For example, for the argument role $a = \text{ENTITY}$ of entity $e = \text{injured man}$,

$$\begin{aligned} C(\langle a, e \rangle, o) &= c(\text{ENTITY of TRANSPORT}, \text{img}) \\ &+ c(\text{injured man}, \text{img}) + c(\text{PERSON}, \text{PERSON}). \end{aligned}$$

The optimal $\mathbf{T} \in \mathbb{R}_+^{n \times m}$ that solves $d(G_t, G_i) = \min_{\mathbf{T}} \mathbf{T} \odot \mathbf{C}$ can be approximated by a differentiable Sinkhorn-Knopp algorithm [5, 29] following [35],

$$\mathbf{T} = \text{diag}(\mathbf{p}) \exp(-\mathbf{C}/\gamma) \text{diag}(\mathbf{q}),$$

where $\mathbf{p} \in \mathbb{R}_+^{n \times 1}$ and $\mathbf{q} \in \mathbb{R}_+^{m \times 1}$. Starting with any positive vector \mathbf{q}^0 to perform the following iteration:

$$\begin{aligned} \text{for } i = 0, 1, 2, \dots \text{ until convergence,} \\ \mathbf{p}^{i+1} = \mathbf{1} \oslash (\mathbf{K} \mathbf{q}^i), \quad \mathbf{q}^{i+1} = \mathbf{1} \oslash (\mathbf{K}^\top \mathbf{p}^{i+1}), \end{aligned}$$

where \oslash denotes element-wise division. $\mathbf{K} = \exp(-\mathbf{C}/\gamma)$. A computational \mathbf{T}^k can be obtained by iterating for a finite number k times,

$$\mathbf{T}^k := \text{diag}(\mathbf{p}^k) \mathbf{K} \text{diag}(\mathbf{q}^k).$$

2.4. Contrastive Learning Objective

We optimize the cosine similarity between image i and positive description t^+ to be close to 1, while negative descriptions t^- to be close to 0,

$$L_1 = \sum_{\langle t, i \rangle} D_{KL}(s(t, i) \parallel \mathbb{1}_{t \in T^+}),$$

where $D_{KL}(\cdot \parallel \cdot)$ is the Kullback-Leibler divergence, and $\mathbb{1}_{t \in T^+}$ is the indicator function showing whether the description is a positive description. It enables our model to handle any number of positive and negative descriptions. Also, we include the descriptions of other images in the same batch as negative descriptions.

We also minimize the distance between two event graphs,

$$L_2 = \sum_{\langle t, i \rangle} d(G_t, G_i).$$

The contrastive learning of event and argument description and the alignment of event graphs are jointly optimized:

$$L = \lambda_1 L_1 + \lambda_2 L_2.$$

We set λ_1 and λ_2 as 1 in this paper.

3. Evaluation Tasks

3.1. Multimedia Event Extraction (M²E²)

Task Setting. Multimedia Event Extraction [17] aims to (1) classify images into eight event types, and (2) localize argument roles as bounding boxes in images. We choose this task as a direct assessment of event structure understanding.

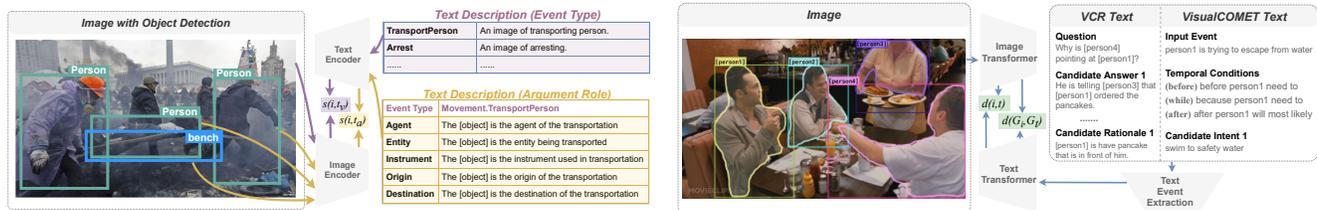
Our Approach. Zero-shot Setting: We perform zero-shot evaluation to show the models' ability to handle open vocabulary events, as required by real-world applications. Also, zero-shot evaluation provides a direct comparison of the effectiveness of event knowledge encoding during pretraining. As shown in Fig. 4a, we select the event type as the one having the highest similarity score $s(i, t)$ with the image, and for each bounding box, we rank candidate argument roles of the selected event type. Supervised Setting: We include the supervised setting to prove the effectiveness of the model architecture at encoding event knowledge in the presence of direct supervision. We use the same training dataset SWiG [25] with 125k images to further finetune our model to compare with the supervised models. During finetuning, we replace the text event extraction results with the annotated events for images, and set the optimal transport plan as the ground truth alignment between arguments and object bounding boxes.

Evaluation Metrics. We follow [17] to use F-scores to evaluate event typing and argument extraction.

3.2. Grounded Situation Recognition (GSR)

Task Setting. Grounded Situation Recognition [25] selects an event type from 504 verbs, and predicts the entity name and the bounding box for each argument role. It is also a direct evaluation of event structure understanding, but with larger size of event types and argument roles.

Our Approach. Similar to Multimedia Event Extraction, in Fig. 4a, we encode 504 candidate verbs using the Text Transformer, and select the top verb as the predicted event



(a) Architecture of event extraction (M^2E^2 and GSR). Event typing (in purple) ranks event type descriptions given the image, and argument extraction (in yellow) rank argument descriptions given the bounding box.

(b) Architecture of VCR and VisualCOMET. We rank ⟨question, answer⟩ and ⟨input event, temporal condition, intent⟩ respectively given the image. We calculate both the image-text level alignment and the event graph alignment.

Figure 4. Architecture for evaluation tasks.

type. For argument extraction, we employ objects detected in the image to rank argument roles, and obtain the union bounding box of objects playing the same argument role. Also, we add the supervised setting similar to M^2E^2 task.

Evaluation Metrics. We follow [25] to evaluate the accuracy of verb prediction (*verb*), argument name prediction (*value* for each argument and *value-all* for all arguments of an event), and argument bounding box and name prediction (*ground* for each argument and *ground-all* for all arguments).

3.3. Image Retrieval

Task Setting. Image retrieval ranks images for the given caption, which is a direct evaluation on the alignment.

Our Approach. We perform the alignment of image and text $d(i, t)$, and also the alignment of event graphs across two modalities $d(G_i, G_t)$.

Evaluation Metrics. We use conventional image retrieval measures including $Recall@1$, $Recall@5$ and $Recall@10$.

3.4. Visual Commonsense Reasoning (VCR)

Task Setting. Given a question, the task contains two sub-tasks: (1) *Answer Prediction* from four options; (2) *Rationale Prediction* from four options to support the aforementioned answer. We include this task to evaluate whether event understanding can better support downstream tasks. To evaluate the quality of pretraining models, we adopt zero-shot settings solely relying on image-text alignment for a fair comparison.

Our Approach. As shown in Fig. 4b, we compute the similarity score between the answer text t and the image i , using both image alignment $d(i, t)$ and event graph alignment $d(G_i, G_t)$. We also consider the question as query and concatenate them with the answer during ranking.

Evaluation Metrics. We evaluate F-scores for both of answer prediction and rationale prediction following [40].

3.5. Visual Commonsense Reasoning in Time

Task Setting. Given the image and the event happening in the image with its participants, VisualCOMET [24] aims to generate “intents” showing what the participants “*need to do*” before the image event, “*want to do*” during the image event, and “*will most likely to do*” after the image event. It

necessitates a deep grasp of events and their connections, as well as a thorough comprehension of arguments roles.

Our Approach. As shown in Fig. 4b, for each image and participant, we use intents from the training data as candidate intents, and rank them based on both image alignment $d(i, t)$ and event graph alignment $d(G_i, G_t)$. The text is the concatenation of (1) input event description, (2) a temporal description detailed in Fig. 4b, and (3) the candidate intents.

Evaluation Metrics. We adopt $Accuracy@50$ following the perplexity evaluation of the state-of-the-art model [24].

4. Experiments

4.1. Pretraining Details

A New Dataset. We collect 106,875 image-captions that are rich in events from news websites [1]. It provides a new challenging image-retrieval benchmark, where each sentence may contain multiple events with a complicated linguistic structure. The average caption length is 28.3 tokens, compared to 13.4 for Flickr30k and 11.3 for MSCOCO. The data statistics are shown in Tab. 3, with structural event knowledge is extracted automatically following Sec. 2.1.

Dataset	Split	#image	#event	#arg	#ent
VOANews	Train	76,256	84,120	148,262	573,016
	Test	18,310	21,211	39,375	87,671
	No-event	12,309	-	-	-

Table 3. Data statistics of VOANews.

Parameter Settings. We utilize the Text and Vision Transformers of “ViT-B/32” to initialize our encoders. More details are included in Appendix.

4.2. Baselines

State-of-the-art Multimedia Pretraining Models. We compare with CLIP [26] by running the public release of “ViT-B/32” and report its scores in the following experiments for a fair comparison. We further pretrain CLIP using the image-captions in the same dataset in Tab. 3 for a fair comparison in terms of data resources.

Setting	Model	Multimedia Event Extraction (M ² E ²)						Grounded Situation Recognition (SWiG)				
		Event			Argument			Event	Argument			
		P	R	F ₁	P	R	F ₁		verb	value	value-all	ground
Zero-shot	CLIP	29.5	65.7	40.7	9.2	12.7	10.7	28.3	19.3	12.8	13.2	4.1
	CLIP pretrained on news	31.7	64.7	42.6	9.7	13.1	11.1	29.9	20.1	13.2	13.9	5.3
	CLIP-Event	36.4	70.8	48.1	13.9	16.0	14.8	31.4	23.1	14.1	14.8	6.3
	w/o OptimalTransport	35.0	59.3	44.1	11.0	12.6	11.9	30.2	21.6	13.6	14.1	5.4
	Single Template	32.3	71.4	44.4	11.9	15.6	13.2	30.4	22.0	13.5	14.0	5.7
	Composed Template	33.9	72.8	46.3	12.7	15.3	13.9	30.9	22.5	13.9	14.2	5.8
	Continuous Prompt	33.6	75.7	46.5	11.1	16.7	13.3	30.4	21.9	13.0	14.1	5.9
	Caption Editing	30.9	71.4	43.2	11.6	13.8	12.6	30.1	20.5	13.0	13.8	5.3
GPT-3 Prompt	34.2	76.5	47.3	12.1	16.8	14.1	31.1	22.1	13.2	14.2	6.0	
Supervised	State-of-the-Art [17, 25]	43.1	59.2	49.9	14.5	10.1	11.9	39.9	31.4	18.9	24.9	9.7
	CLIP finetuned on SWiG	38.1	71.6	49.8	20.9	12.8	15.9	42.6	32.6	19.2	25.2	10.2
	CLIP-Event^{+SWiG}	41.3	72.8	52.7	21.1	13.1	17.1	45.6	33.1	20.1	26.1	10.6
	w/o OptimalTransport	40.3	71.3	51.5	20.8	13.0	16.0	44.7	32.9	19.4	24.4	10.1

Table 4. Evaluation results and ablation studies on image event extraction. We follow the evaluation measures (%) of each benchmark.

Model	Flickr30k		MSCOCO		VOANews	
CLIP	62.2	81.9	30.3	50.3	21.2	23.4
CLIP pretrained on news	64.3	81.2	32.2	50.8	23.5	25.1
CLIP-Event	67.0	82.6	34.0	51.3	27.5	28.7
w/o OptimalTransport	65.6	80.5	32.5	51.0	25.5	26.9

Table 5. R@1(%) on text-to-image (left) and image-to-text (right) retrieval on Flickr30k (1k test), MSCOCO (5k test) and VOANews.

Model	VCR		VisualCOMET
	Answer F ₁	Rationale F ₁	
Perplexity in [24]	-	-	18.2
CLIP	51.1	46.8	20.1
CLIP pretrained on news	51.8	47.2	20.9
CLIP-Event	52.4	49.2	22.4
w/o OptimalTransport	52.0	48.6	21.1

Table 6. Results (%) on zero-shot VCR and VisualCOMET.

State-of-the-art Event Extraction Models. The state-of-the-art event extraction models, such as WASE [17] for Multimedia Event Extraction task, JSL [25] for Grounded Situation Recognition task.

Ablation Study: CLIP-Event w/o Optimal Transport is included as a variant of our model in which we remove the alignment between event graphs. It is trained only on the contrastive loss L_1 .

Ablation Study: Each Prompt Function is used solely during training, for the purpose of comparing its effectiveness.

4.3. Analysis on Event Extraction Tasks

Under zero-shot settings, we achieve 5.5% absolute F-score gain on event extraction, and 33.3% relative gain on argument extraction on M²E², as shown in Tab. 4.

The gains achieved by pretraining on news data are significantly amplified with the help of structural event knowledge. For example, CLIP pretrained on news achieves 1.9% improvement compared to the vanilla CLIP on M²E². Our CLIP-Event significantly boosts the gain to 3.89 times.

Zero-shot CLIP-Event outperforms the state-of-the-art weakly supervised model on argument extraction on M²E² dataset, showing that the proposed optimal transport alignment effectively captures the argument structures, which previous vision-language pretraining models fail.



(a) An example result on M²E². (b) An example result on SWiG.

Figure 5. Example results of event extraction tasks.

For argument localization, CLIP-Event achieves a higher gain on M²E² than SWiG, due to the fact that SWiG uses a different argument bounding box grounding strategy. SWiG merges all objects that play the same role into a single large

bounding box. As shown in Fig. 5b, our approach detects argument roles for each object first, and then merges those objects of the same role into the a large bounding box. In comparison, M²E² allows multiple objects with the same argument role, which is consistent with our approach to use objects aligning with argument roles, as shown in Fig. 5a.

4.4. Analysis on Downstream Tasks

Image Retrieval. (1) VOANews presents a greater challenge due to the various events in the captions and the more difficult sentence structures compared to Flickr30k and MSCOCO, as shown in Fig. 6. The improvement on VOANews is much higher than the gains on Flickr30k and MSCOCO, proving that our model is capable of handling lengthy sentences, particularly those with many events.

Investigators inspect parts of a destroyed car at the site of a car bombing in Beirut, Jan. 21, 2014.



Figure 6. Example results of text-to-image retrieval on VOANews, with the visualizations of the optimal transport plan.

(2) Downstream tasks benefit from fine-grained event graph alignments. For example, in Fig. 6, the strong alignment between objects and *investigators* and *destroyed car* enables the image to be successfully ranked higher.

VCR. (1) On VCR, the rationale F₁ improves more than answer F₁. Rationale prediction is more challenging since it refers to the details of the scene, which our fine-grained alignment well captures. (2) Event knowledge is particularly beneficial for downstream tasks. In Fig. 7, only the correct answer corresponds to the event type of the input image.

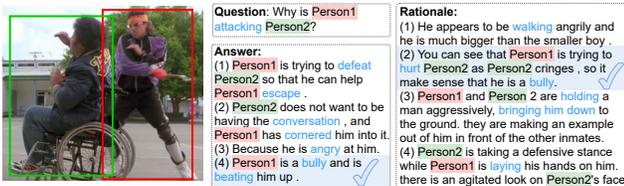


Figure 7. VCR can benefit from event (in blue) understanding.

VisualCOMET. We compare our results to the perplexity of the state-of-the-art model, which is also retrieval-based. The baseline is trained using the training set of VisualCOMET, but our model is an unsupervised model, which achieves superior performance, demonstrating that our model is capable of comprehending events in the images.

4.5. Ablation Studies

Effect of Event Graph Alignment via Optimal Transport. (1) Removing optimal transport (“w/o OptimalTransport”) generally lowers the performance on all evaluation tasks, since it ignores the event graph structures and their cross-media alignment, but relies solely on the overly simplistic image and sentence features. (2) The performance gain on argument extraction task is the highest, since it requires the fine-grained alignment of text and images. (3) We visualize the transport plan in Fig. 6 to bring insights into the learned alignment. It is a global decision that takes the argument structures of two event graphs into account. Thus, distinct argument roles tend to be associated with diverse objects with different visual features in order to achieve a low *global* transport cost. For instance, *investigators* match objects dressed in white, but not soldier objects, due to the dissimilar visual features. Additionally, one argument role tends to be aligned with objects that have similar visual features, e.g., two *investigators* are both dressed in white protection suits.

Comparison between prompt functions. As shown in Tab. 4, GPT3 provides the optimal performance among prompt functions. It leverages the knowledge encoded in GPT3, thus generating natural descriptions with precise event information. Other prompt functions also demonstrate their effectiveness in supporting event understanding.

5. Related Work

Vision-Language Pretraining. Recent years have witnessed great success in Vision-Language pretraining models [4, 11, 12, 14, 18, 22, 26, 31, 34, 41, 42] based on Transformer architectures [32]. Image structures have been proven useful to pretraining models, such as scene graphs [38]. However, event structural knowledge is not well captured in pretraining models, demonstrating deficiencies in tasks related to verb comprehension [10]. We are the first to encode structural event knowledge to enhance vision-language pretraining.

Visual Event Understanding. Previous work simplifies visual events as *verbs* using Subject-Verb-Object triples [2, 6, 9, 13, 19, 21, 28, 30, 33, 36, 43]. Situation Recognition [25, 37] aims to detect argument roles and Multimedia Event Extraction [17] categorizes verbs into event types. However, their limited event ontologies fail to handle open-world events in real applications. In contrast, our proposed pretraining model supports zero-shot event extraction and demonstrate good performance on other downstream tasks requiring image event reasoning.

Cross-media Alignment. Existing pretraining models [3, 4, 18, 31, 41] maximize the alignment across two modalities without taking into account the structure of text and images. Image structures [17, 39] that are analogous to text linguistic structures are proposed. There is, however, a gap between

complicated linguistic structures and image structures. We propose to use the text event graph structures to fill in the gap and compute a global alignment over two event graphs.

6. Conclusions and Future Work

This paper proposes to integrate structural event knowledge into vision-language pretraining. We perform cross-media transfer of event knowledge, by automatically extracting event knowledge from captions and supervising image event structure understanding via contrastive learning. We generate hard negatives by manipulating event structures based on confusion matrices, and design event prompt functions to encode events into natural sentences. To transfer argument structural knowledge, we propose an event graph alignment loss via optimal transport, obtaining a global alignment based on argument structures. It outperforms the state-of-the-art vision-language pretraining models on event extraction and downstream tasks under zero-shot settings. In the future, we will expand this capability to videos to comprehend the evolution of events using argument tracking.

References

- [1] VOA News. <https://www.voanews.com/>. 2, 6
- [2] Yu-Wei Chao, Zhan Wang, Yugeng He, Jiakuan Wang, and Jia Deng. Hico: A benchmark for recognizing human-object interactions in images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1017–1025, 2015. 8
- [3] Liqun Chen, Zhe Gan, Yu Cheng, Linjie Li, Lawrence Carin, and Jingjing Liu. Graph optimal transport for cross-domain alignment. In *International Conference on Machine Learning*, pages 1542–1553. PMLR, 2020. 8
- [4] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020. 1, 8
- [5] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages 2292–2300, 2013. 5
- [6] Pradipto Das, Chenliang Xu, Richard F Doell, and Jason J Corso. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2634–2641, 2013. 8
- [7] Ritendra Datta, Dhiraj Joshi, Jia Li, and James Z Wang. Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys (Csur)*, 40(2):1–60, 2008. 2
- [8] Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30(4):681–694, 2020. 3
- [9] Saurabh Gupta and Jitendra Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015. 8
- [10] Lisa Anne Hendricks and Aida Nematzadeh. Probing image-language transformers for verb understanding. *arXiv preprint arXiv:2106.09141*, 2021. 1, 8
- [11] Zhicheng Huang, Zhaoyang Zeng, Yupan Huang, Bei Liu, Dongmei Fu, and Jianlong Fu. Seeing out of the box: End-to-end pre-training for vision-language representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12976–12985, 2021. 8
- [12] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *arXiv preprint arXiv:2102.05918*, 2021. 1, 8
- [13] Keizo Kato, Yin Li, and Abhinav Gupta. Compositional learning for human object interaction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 234–251, 2018. 1, 8
- [14] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *ICML*, 2021. 8
- [15] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4. *International Journal of Computer Vision*, 128(7):1956–1981, 2020. 3
- [16] Manling Li, Alireza Zareian, Ying Lin, Xiaoman Pan, Spencer Whitehead, Brian Chen, Bo Wu, Heng Ji, Shih-Fu Chang, Clare Voss, et al. Gaia: A fine-grained multimedia knowledge extraction system. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 77–86, 2020. 3
- [17] Manling Li, Alireza Zareian, Qi Zeng, Spencer Whitehead, Di Lu, Heng Ji, and Shih-Fu Chang. Cross-media structured common space for multimedia event extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2557–2568, 2020. 2, 5, 7, 8
- [18] Xiujuan Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020. 1, 8
- [19] Yong-Lu Li, Siyuan Zhou, Xijie Huang, Liang Xu, Ze Ma, Hao-Shu Fang, Yanfeng Wang, and Cewu Lu. Transferable interactivity knowledge for human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3585–3594, 2019. 1, 8
- [20] Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. A joint neural model for information extraction with global features. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7999–8009, 2020. 3
- [21] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection with language priors. In

- European conference on computer vision*, pages 852–869. Springer, 2016. 8
- [22] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *arXiv preprint arXiv:1908.02265*, 2019. 8
- [23] Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60, 2014. 3
- [24] Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. Visualcomet: Reasoning about the dynamic context of a still image. In *European Conference on Computer Vision*, pages 508–524. Springer, 2020. 2, 6, 7
- [25] Sarah Pratt, Mark Yatskar, Luca Weihs, Ali Farhadi, and Aniruddha Kembhavi. Grounded situation recognition. In *European Conference on Computer Vision*, pages 314–332. Springer, 2020. 2, 5, 6, 7, 8
- [26] 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. *arXiv preprint arXiv:2103.00020*, 2021. 1, 2, 3, 4, 6, 8
- [27] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 3
- [28] Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. Foil it! find one mismatch between image and language caption. *arXiv preprint arXiv:1705.01359*, 2017. 8
- [29] Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices. *The annals of mathematical statistics*, 35(2):876–879, 1964. 4, 5
- [30] Jaeyong Sung, Colin Ponce, Bart Selman, and Ashutosh Saxena. Unstructured human activity detection from rgb-d images. In *2012 IEEE international conference on robotics and automation*, pages 842–849. IEEE, 2012. 1, 8
- [31] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019. 1, 8
- [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 8
- [33] Tiancai Wang, Tong Yang, Martin Danelljan, Fahad Shahbaz Khan, Xiangyu Zhang, and Jian Sun. Learning human-object interaction detection using interaction points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4116–4125, 2020. 1, 8
- [34] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904*, 2021. 8
- [35] Hongteng Xu, Dixin Luo, and Lawrence Carin. Scalable gromov-wasserstein learning for graph partitioning and matching. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 3046–3056, 2019. 5
- [36] Bangpeng Yao and Li Fei-Fei. Modeling mutual context of object and human pose in human-object interaction activities. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 17–24. IEEE, 2010. 1, 8
- [37] Mark Yatskar, Luke Zettlemoyer, and Ali Farhadi. Situation recognition: Visual semantic role labeling for image understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5534–5542, 2016. 8
- [38] Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. Ernie-vil: Knowledge enhanced vision-language representations through scene graph. *arXiv preprint arXiv:2006.16934*, 1:12, 2020. 8
- [39] Alireza Zareian, Svebor Karaman, and Shih-Fu Chang. Weakly supervised visual semantic parsing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3736–3745, 2020. 8
- [40] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6720–6731, 2019. 2, 6
- [41] Pengchuan Zhang, Xiujuan Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2021. 1, 8
- [42] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13041–13049, 2020. 8
- [43] Tianfei Zhou, Wenguan Wang, Siyuan Qi, Haibin Ling, and Jianbing Shen. Cascaded human-object interaction recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4263–4272, 2020. 1, 8