

# Critical Analysis: CLIP-Event: Connecting Text and Images with Event Structures

Kai-Yu Lu

2025/11/2

## 1 Methodological Strengths

### 1.1 Event-centric reformulation of vision–language pretraining

The paper introduces an explicit shift from entity-centric to event-centric vision–language pretraining by modeling event types and argument roles as first-class supervision signals. The formulation treats an image and its caption not only as an unstructured pair, but as two parallel realizations of an underlying event graph with nodes for triggers and arguments and edges for semantic roles. This framing is well aligned with downstream tasks that require identifying who did what to whom, where, and with what instrument, such as multimedia event extraction and grounded situation recognition.

The use of a fine-grained event ontology with 187 event types and structured argument inventories provides a rich label space that goes beyond coarse verb labels. This design allows the pretrained model to distinguish visually similar scenes that correspond to different event semantics or different role assignments for the same entities, such as protesters as aggressors versus victims.

### 1.2 Distant supervision via text information extraction

A central methodological strength lies in the use of text information extraction (IE) as a source of distant supervision. Instead of manually annotating event structures in images, the study leverages a state-of-the-art IE pipeline to extract event triggers, types, arguments, and entity coreference chains from captions. These structures are then treated as noisy but large-scale supervision for visual event understanding.

This strategy has several advantages. First, it scales to more than one hundred thousand image–caption pairs without requiring image-side manual event annotations. Second, it aligns with the nature of news data, where captions often describe the same real-world events as the associated images. Third, it reuses mature IE technology that already achieves strong performance on textual event extraction benchmarks, which raises the quality of the weak labels compared with naive heuristic extraction.

### 1.3 Construction of an event-rich pretraining corpus

The VOANews corpus constructed in the paper is methodologically valuable. It contains 106,875 news images with long captions (average length 28.3 tokens) and dense event structures, with more than 80,000 events and nearly 150,000 argument mentions in the training split alone. Compared with MSCOCO and Flickr30k, whose captions are short and object-focused, VOANews provides a challenging setting where captions often contain multiple events, nested clauses, and fine-grained argument details.

By releasing both the images and the automatically extracted event structures, the study creates a resource that supports both pretraining and rigorous evaluation. The statistics reported for event counts, argument counts, and entity counts demonstrate transparency about the scale and richness of the data.

#### 1.4 Event-structure-driven negative sampling and prompting

Another methodological strength is the design of event-structure-driven negative examples. The negative sampling is informed by confusion matrices derived from CLIP predictions and from a text argument extractor. This yields hard negatives that are visually and semantically close to the gold events, such as confusing TRANSPORT with ARREST or rotating argument roles within the same event.

The use of multiple prompt functions to translate event structures into text descriptions is also methodologically robust. The paper explores:

- Single-sentence templates that encode all roles in one description.
- Composed templates that break roles into short declarative sentences.
- Continuous prompts that insert learnable tokens as structural markers.
- Caption editing that minimally perturbs original captions to encode positive or corrupted event structures.
- GPT-3-based prompts that generate fluent descriptions conditioned on structured inputs.

This diversity of prompts reduces overfitting to any single textual pattern and encourages the model to learn event semantics rather than template idiosyncrasies. The ablation study that trains with each prompt function separately provides evidence that every prompt family contributes meaningful signal.

#### 1.5 Multi-level graph alignment via optimal transport

The introduction of event graph alignment based on optimal transport is a central methodological contribution. The model defines cost functions for aligning event nodes and argument nodes across text and image graphs, and uses Sinkhorn-regularized optimal transport to compute a soft alignment plan. This approach captures global structural consistency, since the transport plan is optimized over all possible node correspondences under a shared cost matrix.

Entity-level distances incorporate both mention embeddings and type embeddings, while argument-level costs combine role descriptions, entity mentions, and object features. This multi-component cost design enforces that the same visual object must simultaneously match the textual entity and its role description to achieve low cost, which directly supports fine-grained argument grounding.

#### 1.6 Joint loss combining instance and structure levels

The training objective combines description-level contrastive learning and graph-level alignment. The use of a Kullback–Leibler divergence loss to push positive description similarities toward one and negative similarities toward zero allows varying numbers of positives and negatives per image and integrates batch negatives naturally. The additional loss on graph-level distances ensures that the representations respect structured event correspondences instead of only global image–text similarity.

The balance between the two losses is kept simple by setting scalar weights to one, which avoids hyperparameter overfitting and demonstrates that the architecture, rather than fine-tuned weighting, drives the gains.

## 1.7 Comprehensive evaluation and ablation design

The experimental design covers several complementary tasks: multimedia event extraction (M2E2), grounded situation recognition (SWiG), image retrieval on three datasets, visual commonsense reasoning (VCR), and temporal intent reasoning (VisualCOMET). The study evaluates both zero-shot and supervised settings for event extraction and systematically compares CLIP, CLIP pretrained on the same VOANews data, CLIP-Event, and an ablation without optimal transport.

The ablation suite includes variants trained with single prompt functions and a version without graph alignment. This structure allows clear attribution of performance gains to event structural supervision, graph-level alignment, and prompt diversity. The reporting of precision, recall, and F1 for event and argument extraction, as well as Recall@K and Accuracy@50 for retrieval and reasoning tasks, provides a well-rounded view of performance.

## 2 Key Limitations

### 2.1 Dependence on automatic text information extraction

A major limitation is the heavy reliance on a single text IE pipeline as the source of event supervision. The appendix reports that the event trigger F1 is around the low seventies and argument F1 is substantially lower for some benchmarks. Errors in triggers, types, and roles propagate directly into the pseudo-labels used for pretraining, and there is no explicit mechanism to model or mitigate this noise. The model therefore learns from a mixture of correct and incorrect event graphs.

Because the text IE system is tuned to a specific ontology, the extracted structures are constrained by that design. Events not well covered by the ontology, or non-standard linguistic realizations, are either misclassified or ignored. This introduces systematic bias in the supervision signal and limits the coverage of event types.

### 2.2 Domain and ontology constraints

The VOANews corpus is drawn from newswire style data and focuses on visually salient events such as attacks, protests, disasters, arrests, and transport. Events that are abstract, mental, or conversational are deliberately filtered out in the construction process to avoid non-visual supervision. As a result, the pretrained representation is specialized to a narrow distribution of news events and does not directly address domains such as social media, medical imagery, or scientific diagrams.

The event ontology itself is fixed and fine-grained within the DARPA AIDA schema, which reflects a particular view of geopolitical and disaster events. This yields high coverage for the chosen domain but limits transfer to event schemas in other communities or emerging event types that do not fit the predefined hierarchy.

### 2.3 Visual grounding limitations from object detection

On the visual side, the model depends on Faster R-CNN trained on Open Images for object detection and uses bounding boxes as the basic units of argument grounding. This setup inherits all the limitations of the detector, including missed detections for small, occluded, or rare objects, and label mismatches between Open Images categories and event argument expectations.

Furthermore, roles such as PLACE or abstract roles such as CRIME are difficult to ground in bounding boxes. They are either approximated by coarse regions or omitted. This limits the completeness of visual event graphs and reduces the achievable ceiling on argument alignment metrics.

## 2.4 Limited depth of error analysis and failure characterization

The paper provides qualitative examples of improved alignments and several illustrative cases on M2E2, SWiG, VCR, and VisualCOMET. However, systematic error analysis is limited. The study does not present a detailed breakdown of which event types or argument roles remain difficult, nor does it characterize failure modes such as systematic confusions between visually similar roles or misalignment under cluttered backgrounds.

Without such analysis, it is difficult to assess whether the model struggles primarily with certain ontological categories, specific visual conditions, or particular linguistic constructions. This constrains the diagnostic insight that can be drawn from the reported aggregate metrics.

## 2.5 Computational cost and scalability considerations

The optimal transport based graph alignment introduces non-trivial computational overhead. Even with entropic regularization and a finite number of Sinkhorn iterations, the cost scales with the product of the number of nodes in the text graph and the number of nodes in the image graph. In scenes with many detected objects or captions with multiple events and arguments, this can become expensive.

The appendix reports training with ViT-B/32 backbones on multiple GPUs, but the paper does not quantify the additional cost relative to standard CLIP training or analyze the trade-off between alignment accuracy and computational budget. This limits the practical guidance for scaling the approach to larger backbones or larger corpora.

## 2.6 Evaluation scope and metric coverage

Although the evaluation suite is broad, several aspects of performance remain unexplored. The study focuses on F1, Recall@K, and Accuracy@50, which are standard but do not directly capture structural correctness of full event graphs or human-perceived interpretability of alignments. There is no measure of graph edit distance to gold structures or of calibration of similarity scores.

In addition, all downstream tasks are benchmark-centric and there is no user study or application-driven evaluation that measures the impact of event-aware pretraining on actual multimedia analytics workflows, such as news monitoring or misinformation detection.

## 3 Technical Bottlenecks

### 3.1 Complexity of cross-modal graph alignment

The reliance on optimal transport for event graph alignment constitutes a fundamental technical bottleneck. The cost matrix must compare every text node with every image node, and the Sinkhorn iterations must repeatedly normalize the resulting matrix. While manageable for small graphs, this approach becomes increasingly expensive as the number of objects and arguments grows.

This complexity discourages the incorporation of richer scene graphs, such as those including relations between objects or higher-order structures, since the cost would then scale with even larger graph sizes. It also complicates extension to video, where events evolve across time and graphs would involve multiple frames.

### 3.2 Information bottlenecks in object and role representations

The representation of arguments relies on average-pooled features over bounding boxes and contextualized token embeddings over entity mentions. This design introduces information bottlenecks at two levels.

First, average pooling over patches within a bounding box discards spatial structure and fine-grained visual cues that could be informative for role distinctions, such as whether an object is being held, worn, or targeted. Second, encoding roles via short textual descriptions or learned tokens compresses potentially rich semantic information about argument constraints into low-dimensional vectors that may not fully capture the ontology’s nuances.

These bottlenecks may contribute to residual confusion between roles that are visually similar but semantically distinct, such as AGENT versus ENTITY or INSTRUMENT versus PLACE, particularly in crowded scenes.

### 3.3 Integration challenges between IE pipeline and pretrained encoders

The overall pipeline integrates an external IE system, a CLIP-style encoder, and an object detector. Each component has been trained in isolation on different datasets with different objectives. The integration relies on fixed outputs from the IE and detector systems, which are not updated during pretraining.

This separation creates several integration challenges. Misalignments in tokenization or sentence segmentation between IE and the text encoder can lead to suboptimal contextualization of event triggers and entities. Inconsistent naming conventions between the event ontology and object labels can reduce the effectiveness of type-based distance terms. There is no joint training that could adapt the IE representations to better align with the CLIP embedding space.

### 3.4 Trade-offs in negative sampling and prompt diversity

The negative sampling mechanism uses confusion matrices derived from existing models, which biases the negative pool toward errors characteristic of those models. While this yields hard negatives, it also risks reinforcing certain blind spots. For instance, if the base classifier rarely confuses two event types that are conceptually close but underrepresented in the data, the negative sampler will not expose that distinction.

Prompt diversity introduces another trade-off. Template-based prompts are structurally clear but may produce unnatural language, while GPT-3-based prompts are fluent but harder to control and depend on external language modeling capabilities. Maintaining consistency across these prompt families requires careful engineering, and differences in language style may introduce variability in the supervised signal that does not always reflect intended event distinctions.

### 3.5 Limited temporal and discourse modeling

Although the paper evaluates on VisualCOMET, which is a temporal reasoning benchmark, the underlying model remains a static image–text encoder without explicit temporal or discourse modeling. Event relations such as before, after, or causality are encoded only implicitly through the textual descriptions of intents and temporal phrases.

This creates a bottleneck for extending the approach to multi-event narratives or to video, where temporal dynamics and argument trajectories are central. Without explicit temporal state representations or argument tracking mechanisms, the model relies on coarse statistical associations between visual cues and temporal phrases, which may be insufficient for complex scenarios.

## 4 Research Implications

### 4.1 Implications for vision–language pretraining

The empirical gains of CLIP-Event over CLIP and CLIP pretrained on the same data demonstrate that injecting structured event supervision into vision–language pretraining yields substantial benefits across

diverse tasks. This provides concrete evidence that large-scale contrastive learning alone does not fully capture event semantics, and that structured knowledge derived from text can significantly enrich the learned representations.

The results also suggest that modern VLP models should incorporate event ontologies and argument structures as first-class components, analogous to the way some models incorporate scene graphs. Event graphs provide a powerful inductive bias for understanding interactions and roles, which are central to many high-level vision tasks.

## 4.2 Implications for multimedia event extraction and situation recognition

The improvements in zero-shot event extraction and grounded situation recognition indicate that pretraining with event-centric supervision can reduce dependence on task-specific labeled data. This has direct implications for multimedia event extraction systems that must operate in open-world settings where new event types and role configurations continually arise.

The ability to generalize to unseen events via natural language descriptions of types and roles highlights the potential of combining structured ontologies with natural language prompts. This suggests a trajectory toward hybrid systems that unify symbolic event schemas and neural representation learning.

## 4.3 Implications for commonsense and temporal reasoning

The gains on VCR rationale selection and VisualCOMET intent prediction demonstrate that improved event understanding contributes to stronger visual commonsense reasoning. By encoding who does what to whom with which instrument and where, the model acquires a foundation for inferring plausible intentions and outcomes.

This connection underscores that progress in visual commonsense reasoning does not only depend on larger backbone models, but also on better structured supervision grounded in events and roles. Event-aware pretraining can complement textual commonsense resources and support more faithful, role-aware explanations in multimodal tasks.

## 4.4 Implications for benchmark design and evaluation

The construction of VOANews and the observation that improvement is larger on this dataset than on Flickr30k or MSCOCO highlight an important gap between standard captioning datasets and real-world multimedia content. News captions are longer, more syntactically complex, and more event-dense than traditional benchmarks, which often focus on simple object descriptions.

This suggests that future benchmarks for vision–language understanding should better reflect event-centric, multi-sentence narratives with multiple actors and overlapping events. It also supports the view that evaluation should consider structural aspects of understanding, not only coarse image–sentence similarity.

# 5 Potential Research Directions

## 5.1 Joint learning of event extraction and vision–language representations

One natural extension is to replace the fixed IE pipeline with a trainable event extraction module that is jointly optimized with the vision–language encoder. Such a module could share parameters with the text encoder and benefit from cross-modal feedback, potentially improving extraction accuracy in domains where image evidence supports event disambiguation.

Uncertainty-aware training could downweight noisy extractions and allow the model to learn which IE predictions are reliable in which contexts. This direction would reduce dependence on a single frozen IE system and create a more cohesive architecture.

## 5.2 Toward open-schema and schema-light event representations

The current approach is tightly bound to a specific event ontology. Future work can explore schema-light or open-schema event representations, in which event types and roles are represented by textual descriptions without fixed enumerations. This would allow the model to adapt to new event categories and role variants that arise in new domains or languages.

Techniques such as representation learning over event type definitions, clustering of free-text event descriptions, and meta-learning over ontologies could help bridge the gap between fixed schemas and open-world event semantics.

## 5.3 Efficient and scalable graph alignment mechanisms

Given the computational cost of optimal transport, research into more efficient graph alignment mechanisms is warranted. Possible directions include:

- Approximating optimal transport with low-rank or sparse transport plans.
- Using graph neural networks to perform iterative message passing between text and image nodes, learning alignments implicitly rather than solving explicit transport problems.
- Hierarchical alignment strategies that first align coarse structures, such as event types, and then refine alignments at the argument level.

These approaches could preserve the benefits of structural alignment while improving scalability to larger graphs and video sequences.

## 5.4 Richer visual representations for argument roles

To alleviate information bottlenecks, future models can incorporate more expressive visual encoders for arguments, such as region-level transformers that preserve spatial relationships, pose cues, and interaction patterns. Integrating human pose estimation, object affordance models, or relational features between bounding boxes could help distinguish roles that are visually similar at the object category level but differ in configuration.

Combining such rich visual features with structured role descriptions has the potential to improve fine-grained role disambiguation, particularly in complex scenes with multiple similar entities.

## 5.5 Temporal and multi-event modeling

Extending CLIP-Event to video and multi-image narratives is a promising direction. This would require event and argument tracking across time, modeling of event sequences, and representation of temporal relations such as cause, enablement, and prevention.

Temporal event graphs with edges representing before/after and causal links could be aligned across modalities, and contrastive learning could be extended to sequences of events rather than single events. This line of work would directly benefit tasks such as video event detection, story understanding, and temporal commonsense reasoning.

## 5.6 New evaluation methodologies and interpretability tools

Finally, there is a need for evaluation methodologies that specifically target structural correctness and interpretability of event understanding. Potential directions include:

- Graph-level metrics that assess correctness of full event graphs, not only individual roles.
- Human evaluations that rate the faithfulness of textual rationales with respect to visual event structures.
- Visualization tools for transport plans and alignments that allow practitioners to inspect model behavior and diagnose systematic errors.

Such tools would support both scientific understanding of model capabilities and practical deployment in settings where interpretability is critical.

## 6 Conclusion

The CLIP-Event framework provides a compelling demonstration that integrating structural event knowledge into vision–language pretraining yields substantial gains in multimedia event extraction, situation recognition, image retrieval, and visual commonsense reasoning. The methodological strengths include the exploitation of text-based event extraction for distant supervision, the construction of an event-rich news dataset, the design of event-structure-driven negative sampling and prompt families, and the use of optimal transport for cross-modal event graph alignment.

At the same time, the approach exhibits clear limitations, including dependence on a fixed and imperfect IE pipeline, domain and ontology constraints, computational costs associated with graph alignment, and limited depth of error analysis. These limitations point to technical bottlenecks in integrating heterogeneous components, scaling structural alignment, and generalizing beyond a specific event schema and domain.

The most promising research directions include joint learning of event extraction and vision–language representations, development of open-schema event models, design of more efficient graph alignment mechanisms, enhancement of visual representations for argument roles, explicit temporal modeling, and richer evaluation protocols for structural understanding and interpretability. Pursuing these directions can advance the broader goal of building multimodal systems that reason over events and participants in a manner that is both robust and semantically faithful.