# Filming events[1]

Emar MAIER — *University of Groningen*

**Abstract.** Eckardt argues against the ontological reduction of events to "little movies in time and space." In this paper I explore what this means for the representation of events in visual discourse, specifically film. As it turns out, we can build a rather intuitive film semantics on top of the 'regional event' ontology that Eckardt rejects. But we can also follow Eckardt's reasoning and incorporate her participant-based event ontology.

**Keywords:** events, pictorial meaning, film semantics, discourse, narration

## 1. Introduction

In her book *Adverbs, events, and other things*, Eckardt (1998) investigates what kind of event ontology is needed to do natural language semantics. After a swift rejection of the reduction of events to temporal intervals, she considers a prima facie more sensible reduction of events to spatio-temporal regions, i.e., events as "little movies in time and space" (Eckardt, 1998). She then rejects this view on the grounds that different events can occur simultaneously, occupying the exact same place, like A selling a book to B while B is buying that book from A.

Today, events are considered essential ingredients in models for natural language interpretation. We need them in our logical forms to capture the interpretation of adverbs, tense, aspect, speech reports, etc. Beyond the sentence, they play an equally essential role. Narrative, for instance, is traditionally defined as the representation of a sequence of events (Labov, 1972). More generally, in discourse coherence theories like SDRT, many discourse relations (NARRATION, RESULT, ELABORATION, BACKGROUND, etc.) are interpreted as establishing relations between events, and segmentation of text into elementary discourse units is partly driven by the assumption that such units should introduce and describe a main eventuality (Asher and Lascarides, 2003).

Now that semantics, and especially discourse semantics, is increasingly applied to multimodal and visual sign systems (Patel-Grosz et al., 2023), semanticists need to consider event representation in non-verbal signs and media. This super-linguistic extension of (discourse) semantics is not trivial. In verbal discourse segmentation, for instance, we take clauses as basic event descriptions, with verbs as a conventional lexical source of event introduction, often modified in various intricate ways by tense, aspect, and adverbs further up in the semantic composition of the clause. In the pictorial domain there is no (obvious) analogue of semantic composition, let alone verbs or aspectual modifiers.

Abusch (2014) and Schlöder and Altshuler (2023) further problematize the application of linguistic discourse theories to comics, i.e., pictures put in sequence to tell coherent stories. Since,

---

they argue, individual pictures depict states of the world, not events, we cannot rely on standard event-based discourse relations like NARRATION to infer narrative progression. Maier (2024a) and LaRose (2024) take a critical look at their argument and conclusion, but in this short paper I want to sidestep the alleged aspectual limitations of static pictures by looking at the representation of events in a different, inherently dynamic pictorial medium, viz. film. Can we use the view of events as regions of space-time to build a proper event-based semantics of film? And what does Eckardt's rejection of regional event ontology mean for this endeavor?

## 2. Film as discourse

Film is a sequence of shots, stitched together in editing to convey a coherent narrative. The discourse semantics view of film interpretation (Cumming et al., 2017; Wildfeuer, 2014) takes these shots as elementary discourse units, i.e., discrete chunks of information analogous to clauses in verbal discourse or panels in comics. Interpreting a film as a coherent discourse means that the interpreter connects the shots via implicit discourse relations, like NARRATION, BACKGROUND, or ATTRIBUTION.

Consider the following 3 consecutive shots from a scene in *Basic Instinct*, with each shot described in an English sentence.[2]
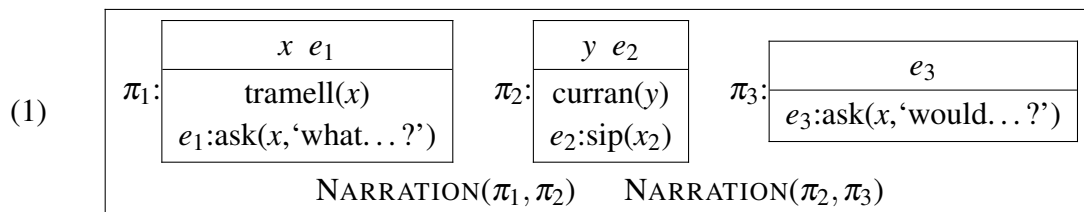


Tramell asks, "What difference does it make?"



Curran takes a sip of his coffee.



Tramell asks him, "Would you like a cigarette, Nick?"

If we look just at my linguistic retelling and analyze its discourse structure according to standard SDRT parsing algorithms (Asher and Lascarides, 2003), we'd get (at least) three elementary units ($\pi_1$, $\pi_2$, $\pi_3$), corresponding to the three sentences, connected by NARRATION relations.[3] In a simplified, boxy (S)DRS representation format:

$$
\text{(1)} \quad \pi_1: \boxed{\begin{array}{c} x\ e_1 \\ \hline \text{tramell}(x) \\ e_1:\text{ask}(x, \text{'what}\ldots?\text{'}) \end{array}} \quad \pi_2: \boxed{\begin{array}{c} y\ e_2 \\ \hline \text{curran}(y) \\ e_2:\text{sip}(x_2) \end{array}} \quad \pi_3: \boxed{\begin{array}{c} e_3 \\ \hline e_3:\text{ask}(x, \text{'would}\ldots?\text{'}) \end{array}}
$$
$$
\text{NARRATION}(\pi_1, \pi_2) \quad \text{NARRATION}(\pi_2, \pi_3)
$$

The semantics of NARRATION($\pi_1, \pi_2$) in (1) says that the primary eventuality introduced by the first unit precedes that introduced by the second (in close spatio-temporal proximity, some more details follow in Section 4.2). Thus, (1) represents the intuitively plausible reading of the linguistic discourse where it describes a sequence of four consecutive events occurring one after the other.

---

[2]Paul Verhoeven, 1992, *Basic Instinct*. https://youtu.be/rAzbU8hayfw?feature=shared&t=229

[3]The direct quotations could be analyzed as separate, subordinated discourse segments, connected via ATTRIBUTION, but we'll ignore that for now (Maier, 2023).

The natural starting point for an SDRT-style discourse analysis of the film sequence is that it has exactly the same structure, i.e., each of these shots constitutes an elementary discourse unit, connected to the next via NARRATION. Interpreting NARRATION then presupposes that the individual shots must not only express propositional contents (by depicting what the world is like), but more specifically that each shot depicts (or otherwise semantically makes available in the discourse representation) a primary event, just like the DRS boxes in (1) do.

Luckily, Abusch's hypothesis about pictures in comics being inherently stative doesn't seem to apply, since unlike panels, shots are themselves dynamic and hence may depict movement and change. But does that mean shots literally depict events? And how many, which ones, and how do we determine the primary event expressed by a shot?

## 3. Events in film semantics

The spatio-temporal view of events, combined with a generalization of the geometric projection view of picture semantics, can help us make sense of filmic event depiction.

### 3.1. Geometric shot semantics

Following Greenberg (2013) and Abusch (2020), the geometric projection view of picture semantics holds that the meaning of a picture is the set of viewpoint-centered worlds $\langle w, v \rangle$ such that $w$ as seen from $v$ looks like that picture. More precisely, $w$ as seen from $v$ looks like picture $p$ iff applying the contextually relevant projection function $\Pi$ to $w$ and $v$ yields $p$. A projection function is a certain type of well-behaved structure-preserving mapping from 3D space (a location in $w$ seen from $v$) to a 2D plane, and it might corresponds to, say, the technique of black-and-white linear perspective line-drawing (Abusch, 2020; Maier, 2024a; Willats, 1997).

A shot is a moving picture, which we might model as a function $s$ from a sequence of consecutive time points (in discrete time, on account of a film's typically finite framerate) to a set of pictures (frames). Like in pictorial semantics, the meaning of a shot is still the set of viewpoint-centered worlds that look like (i.e., that can be geometrically projected onto) the shot, except now the viewpoint is not a vector (the viewing direction) located in time and space, but such a vector moving through space for a certain amount of time in $w$ – think of a camera that may be panning or moving while filming a shot. Mathematical details will be somewhat tedious but the idea is straightforward: the meaning of a shot is the set of worlds that look like that when observed through the lens of an inferred camera running for the duration of the shot.

(2)    $[\![s]\!] = \{ \langle w, v \rangle \mid$ world $w$ looks like shot $s$ when observed from dynamic viewpoint $v \}$

To make that more precise we can spell out the 'looks like' in terms of geometric projection functions. Say, $w$ is a possible world and $v$ is a dynamic viewpoint, i.e., formally, $v$ is a continuous function from some closed temporal interval ($Dom(v)$, the 'runtime' of the 'camera') to a set of spatial viewpoints (or 'camera standpoints', i.e., vectors that have a direction, a viewing angle, and a location in space).

(3)    world $w$ looks like shot $s$ when observed from dynamic viewpoint $v$, relative to geometric projection function $\Pi^*$, iff $\Pi^*(w, v) = s$

Strictly speaking, the geometric projection functions familiar from Greenberg's and Abusch's picture semantics, are defined only for static images and static viewpoints. We can define a dynamic generalization $\Pi^*$ of a static 3D-to-2D projection function $\Pi$ by viewing a shot as a series of still images. In film-like media, shots have a finite framerate so $s$ will consist of, say, $n$ frames $(s_1, \ldots s_n)$. We can then evenly sample an equal number of moments $(t_1^v, \ldots t_n^v)$ from the camera runtime $Dom(v)$ in the world and define dynamic projection as follows:[4]

(4) $\qquad \Pi^*(w, v) = s$ iff $\Pi(w, v(t_i^v)) = s_i$ for all $1 \leq i \leq n$

## 3.2. Events in viewpoint-visible regions

With (2) and (3) we analyze the meaning of a shot as a set of viewpoint-centered worlds. To get to the events depicted by a shot, note first that a viewpoint-centered world determines a region of space-time in a world, viz. the moving, cone-shaped 'search light' area of the world seen by the moving viewpoint.

Next, recall then from Section 1 that the spatio-temporal reduction of events that Eckardt discusses (and dismisses), views events as just such regions of space-time.

It is thus rather straightforward to capture SDRT's notion of a discourse unit's main event. Just like a verb lexically introduces a set of events (with lexical rules like: $[\![walk]\!] = \lambda e.walk'(e)$), the viewpoint-centered worlds that make up the semantic interpretation of a shot also determine a set of 'regional events', viz. those that fall inside the viewpoint-visible space-time regions.

## 3.3. Summing regional events

There will typically be many intuitively distinct events happening inside a single viewpoint-visible region of a single possible world. This happens even if the viewpoint has a very short runtime and comprises no more than a static close-up of a single person against a plain (or blurred out) background. A single such shot could depict Tramell looking up, fixing her gaze on Curran, looking down at her cigarette, saying something, shifting her posture, tilting her head, blinking etc, all in the space of a few seconds.

The same is not always true for linguistic discourse units, as a simple clause like "Curran takes a sip of his coffee" plausibly introduces only an event of sipping.[5] Unlike the shot, the sentence does not entail that he is gazing intently at something, breathing softly, leaning forward, etc.

Even in language, though, we do find elementary discourse units with quantified subjects ("three detectives entered the room") that seem to describe multiple distinct events. Since SDRT's coherence relations demand a unique "primary eventuality" discourse referent, we could assume some mereological structure on events and sum up these small individual events to create a complex maximal event described by the quantified unit. This strategy could be extended to complex discourse units, and to shots. All the small $v$-visible regional events are then

---

[4]We can extend this to 'infinite framerate film', which would mean that the medium itself literally moves, like, perhaps, some form of Wayang shadow puppet theater?

[5]To be clear, the sentence of course does not pick out a unique event, but just a unique event *type*, viz. sipping. A shot tends to pick out many plausibly distinct events of many plausibly distinct event types.
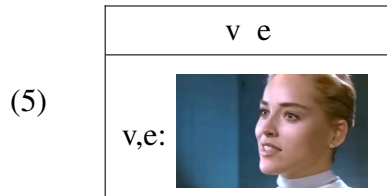
subevents of a single maximal event, which we might call the maximal regional event visible from $v$.[6]

## 4. Interpreting filmic discourse structures

Let's integrate the remarks about filmic event semantics above into a dynamic DRT semantics and then use that to make sense of filmic narration in SDRT.

### 4.1. Shot dynamics

Following Maier and Bimpikou (2019) we start by including pictures into the DRS language in the form of pictorial conditions. In this case we include shots, i.e., moving pictures (but I'll refrain from including animated gifs). We'll introduce discourse referents for moving viewpoints and the maximal events visible from them. A basic DRS for the representation of the first *Basic Instinct* shot above:

(5)



Dynamic semantics of DRS' is often specified as a relation between input and output contexts (where a context is a world–assignment pair, i.e., a possible world enriched with a number of 'discourse referents' that we're keeping track of when updating the common ground with information conveyed by each discourse unit). For a simple moving-picture-DRS like (5) the dynamic semantics extends the domain of the input assignment with discourse referents for a viewpoint, $v$, and an event, $e$, such that input world $w$ observed through $v$ looks like the shot, and $e$ is the maximal $v$-visible regional event in $w$. In (6) I write this more precisely, using the notation $f \subseteq_X g$ to mean that partial assignment $g$ is an extension of $f$ that includes $X$ in its domain (Geurts et al., 2016).

(6)    $\langle w, f \rangle [\![ (5) ]\!] \langle w', f' \rangle$ iff
   a.    $w = w'$ and $f \subseteq_{\{v,e\}} f'$

   b.    $\Pi(w, f'(v)) =$ 

   c.    $f'(e)$ is the maximal regional event visible from $f'(v)$ in $w$.

### 4.2. Narration

In standard SDRT, the coherence relation of NARRATION is a so-called veridical relation. This means that if two units $(\pi_1, \pi_2)$ are connected by NARRATION we have to compute the semantic contributions $(K_{\pi_1}, K_{\pi_2})$ of both those units and add those to the context. We then add the NARRATION-specific part, which says that the main event introduced by $\pi_1$ ($e_{\pi_1}$) precedes that

---

[6]Since regional events are spatio-temporal regions we might even equate this maximal event with the entire $v$-visible region itself.

added by $\pi_2$ (in close spatio-temporal proximity). In other words (with dynamic conjunction $\wedge$) (Asher and Lascarides, 2003):
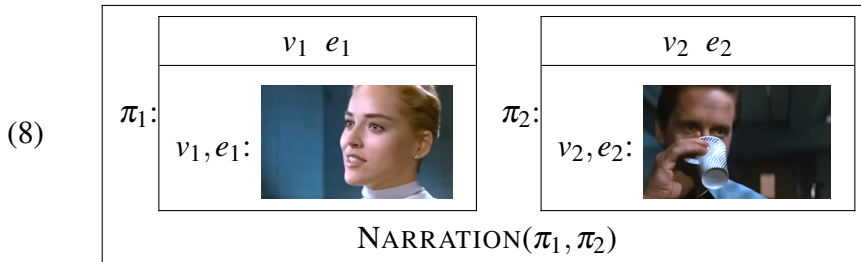
(7)   $[\![\text{NARRATION}(\pi_1, \pi_2)]\!] = [\![K_{\pi_1} \wedge K_{\pi_2} \wedge (e_{\pi_1} \prec e_{\pi_2})]\!]$

The starting point of the discourse semantic view of film interpretation is that discourse relations like NARRATION, inferred between shots in film (or panels in comics), work just like in linguistic discourse, at least semantically. However, there appear to be some additional constraints on the temporal and spatial relations inferred between shots in filmic narrative sequences. For instance, in film (but not in comics or language), temporal progression between shots within a scene is typically assumed to be gapless. We could add this as a further constraint in (7) ($Dom(f(v_{\pi_1})) \supset\subset Dom(f(v_{\pi_2}))$).

Moreover, spatial relations between viewpoints in a narrative scene are further constrained by film conventions like the X-Constraint: the viewpoint of the second shot must be on the same side of a salient action line as that of the first (Cumming et al. 2017). Note that these filmic constraints are relations between viewpoints, which explains, perhaps, why we never encounter them in linguistic discourse. If we stipulate vacuous compliance with viewpoint constraints whenever the discourse units in a narration sequence don't provide viewpoints, we can thus maintain that NARRATION has a completely uniform semantics across modalities.

### 4.3. Example: Filmic narration in *Basic Instinct*

We'll apply the above to the first two shots of the *Basic Instinct* scene from Section 1. We hypothesized a discourse structure with each shot treated as a discourse unit, connected via NARRATION. Here's the box-style representation of the SDRS, with moving image conditions in the DRS boxes:

(8)



The dynamic model-theoretic interpretation of this entire box is simply the interpretation of its only contentful condition, the NARRATION relation, which, as a veridical relation, in turn calls on the interpretations of the two labeled DRS boxes representing the elementary discourse units, as defined in (7). Working out the dynamic semantics gives:

(9)   $\langle w, f \rangle [\![(8)]\!] \langle w', f' \rangle$ iff $\langle w, f \rangle [\![\text{NARRATION}(\pi_1, \pi_2)]\!] \langle w', f' \rangle$ iff
  a.   $w = w'$ and $f \subseteq_{\{v_1, e_1, v_2, e_2\}} f'$
  b.   $\Pi(w, f'(v_1)) =$  and $\Pi(w, f'(v_2)) =$ 
  c.   $f'(e_1)$ is the maximal regional event visible from $f'(v_1)$ and $f'(e_2)$ is the maximal regional event in $w$ that is visible from $f'(v_2)$
  d.   $f'(e_1) \prec f'(e_2)$

e.    (some film-specific viewpoint constraints between $f'(v_1)$ and $f'(v_2)$ such as temporal continuity and the X-Constraint)

We see in (9) that the filmic discourse structure in (8) has the intended interpretation: (i) first the world looks like the first shot, then immediately after that it looks like in the second shot; and (ii) we're adding four discourse referents to the common ground, two for the inferred viewpoints (which allows us to say that they are shot from the same side of the action line given by the first eyegaze direction, which in turn allows us to infer that Tramell and Curran are not looking in the same direction, but (presumably) facing each other), and two for the regional events depicted in the two shots (which allows us to say that these events occur right after each other).

## 5.  Events and their participants

### 5.1.  The limitations of regional events

So far we've been working with regional events, i.e., events as regions of space-time. But as Eckardt points out, many different events can occupy the same space-time region. Take a shot of two people at a store, one buying a pair of socks, the other selling it.[7]

(10)    

The visible region determined by the relevant inferred viewpoint always contains both of these events (or neither), and always many more besides (the customer smiling, Chaplin folding the socks and wrapping them, the woman paying etc). Our strategy of summing them all together into a complex maximal visible event of which all of these distinct events are subevents in a way amounts to giving up on the intuitive idea that elementary discourse units introduce a single primary event.

The maximal event sum strategy doesn't seem to lead to significant problems in simple narration sequences. But then, if you think about it, the events don't really do much work in the semantics of NARRATION in filmic discourse representations either. If we were only interested in narrative progression we could simply do away with events altogether and just stipulate spatial proximity and temporal progression as viewpoint constraints.

Ultimately, we want to use our discourse framework also to capture more interesting coherence relations, like ATTRIBUTION to capture filmic dream sequences, for instance (Maier, 2024b). In order for that to work we want to say that a shot zooming in on a character's sleeping face might introduce a contentful (and, moreover, invisible and hence not obviously regional) event of dreaming (rather than a complex regional event of turning one's head, snoring, breathing, lying in bed, etc).

---

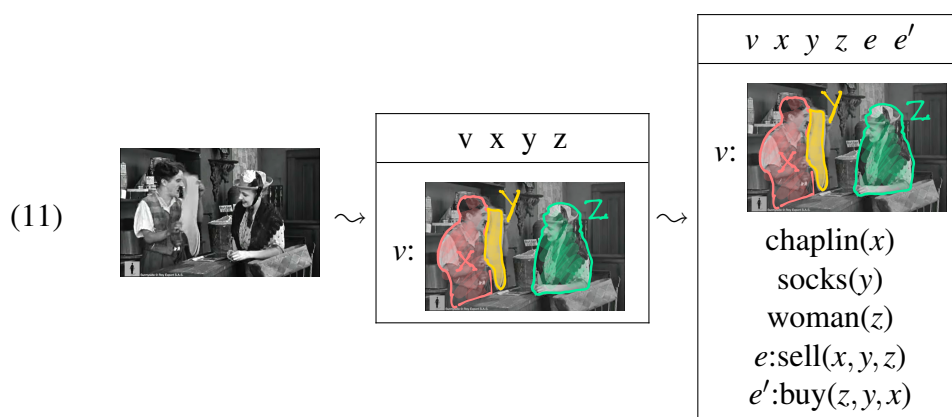[7]Charlie Chaplin, 1919, *Sunnyside* `https://youtu.be/bZwG5r25W1c?si=tRzziBlsEzbZFHyy`

## 5.2. Participant-based individuation

Let's briefly consider Eckardt's proposed solution to the buy/sell problem for regional event ontology. Her idea is that events are ontologically individuated by their semantic participants. There's one event, the selling, in which Chaplin is the agent, the socks are the theme, and the customer is the recipient; and there's an ontologically distinct event, the buying, in which the customer is the agent, socks are the theme, and Chaplin is the source. In terms of visible regions there is no difference, but ontologically there is, because they have different participants in different roles. Building on that idea we might say that we don't literally *see* events; we see regions of space-time and participants and infer events from that.

Following Maier and Bimpikou (2019) our (pre-)semantic processing of images includes a component where our visual system identifies certain salient image regions and labels those with fresh discourse referents. These individual discourse referents are meant to track reference to salient depicted objects and individuals in the discourse and hence can be pragmatically equated across panels.

We can extend this to shots. We identify a finite number of salient moving regions in a shot and label those with fresh discourse referents. These discourse referents can be equated across shots (whenever that helps create a more coherent output representation).

Events are now inferred based on the visible participants and their visible properties and relations to each other. In (10) we can visually identify (at least) three salient regions, corresponding to Chaplin, the socks, and the customer. We don't see a specific region that we can call a buying or a selling event, but we know that the three individuals that we do see are all participating in (i.e., performing, undergoing, experiencing, etc.) a great number of eventualities. Two of the more salient such events are the buying and the selling, so we can add those to the discourse record, the same way we can also add the inferred (but strictly invisible) information that $x$'s name is Charlie Chaplin, and $y$ is a pair of socks.

(11)



In sum, the interpretation of panels and shots proceeds in two stages.

Stage 1: semantic processing:
- image is added to fresh DRS;
- salient image regions are identified;
- fresh discourse referents for the inferred viewpoint and salient regions introduced in DRS universe.

Stage 2: pragmatic processing:
-     infer salient properties, names, relations that feature the discourse referents as arguments.
-     infer events and event-types featuring the discourse referents as participants.

When we venture beyond the single shot, we have to add a third stage of discourse processing:

Stage 3: discourse processing:
-     introduce the DRS as a $\pi_i$-labeled box into the SDRS representing the discourse
-     connect it via suitable coherence relations to previous discourse units
-     equating new discourse referents with accessible discourse referents introduced by previous units.

## 6. Conclusion

Eckardt dismisses the view that events are regions of space-time, on metaphysical/semantic grounds. This is unfortunate because those regional events would be very attractive for super-linguists trying to extend event-based linguistic discourse theories to pictorial media – especially film, since film shots literally depict precisely such regions of space-time.

I've demonstrated how we can indeed build a film discourse semantics on top of a regional event ontology. Following Eckardt's rejection of regional events, I then explored briefly how visual semantics can incorporate her more fine-grained event ontology based on visually identifiable event participants.

## References

Abusch, D. (2014). Temporal succession and aspectual type in visual narrative. In L. Crnič and U. Sauerland (Eds.), *The Art and Craft of Semantics: A Festschrift for Irene Heim*, pp. 9–29. MITWPL.

Abusch, D. (2020). Possible-Worlds Semantics for Pictures. In *The Blackwell Companion to Semantics*. Wiley.

Asher, N. and A. Lascarides (2003). *Logics of Conversation*. Cambridge: Cambridge University Press.

Cumming, S., G. Greenberg, and R. Kelly (2017). Conventions of Viewpoint Coherence in Film. *Philosophers' Imprint 17*(1), 1–29.

Eckardt, R. (1998). *Adverbs, Events, and Other Things: Issues in the Semantics of Manner Adverbs*. Tübingen: Niemeyer.

Geurts, B., D. Beaver, and E. Maier (2016). Discourse Representation Theory. In E. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Spring 2016 ed.). Metaphysics Research Lab, Stanford University.

Greenberg, G. (2013). Beyond Resemblance. *Philosophical Review 122*(2), 215–287.

Labov, W. (1972). *Language in the Inner City: Studies in the Black English Vernacular*. University of Pennsylvania Press. Google-Books-ID: snEEdFKLJ5cC.

LaRose, G. (2024). Revisiting stativity in pictorial narratives. *Proceedings of Sinn und Bedeutung 28*, 526–539.

Maier, E. (2023). Attribution and the discourse structure of reports. *Dialogue & Discourse 14*(1), 34–55.

Maier, E. (2024a). Pictorial language and linguistics. LingBuzz Published In:.

Maier, E. (2024b). Reporting, telling, and showing dreams.

Maier, E. and S. Bimpikou (2019). Shifting perspectives in pictorial narratives. *Sinn und Bedeutung 23*(2), 91–106.

Patel-Grosz, P., S. Mascarenhas, E. Chemla, and P. Schlenker (2023). Super Linguistics: an introduction. *Linguistics and Philosophy 46*(4), 627–692.

Schlöder, J. and D. Altshuler (2023). Super Pragmatics of (linguistic-)pictorial discourse. *Linguistics and Philosophy 46*, 693–746.

Wildfeuer, J. (2014). *Film Discourse Interpretation: Towards a New Paradigm for Multimodal Film Analysis*. Routledge.

Willats, J. (1997). *Art and representation: new principles in the analysis of pictures*. Princeton, N.J.: Princeton University Press.