

Compositionality and Event Retrieval in Complement Coercion: A Study of Language Models in a Low-resource Setting

Anonymous ACL submission

Abstract

In sentences such as *John began the book*, the complement noun, lexically denoting an entity, is interpreted as an event. This phenomenon is known in linguistics as *complement coercion*: the event associated with the verb is not overtly expressed but can be recovered from the meanings of other constituents, context and world knowledge. We investigate whether language models (LMs) can exploit sentence structure and compositional meaning to recover plausible events in complement coercion. For the first time, we tested different LMs in Norwegian, a low-resource language with high syntactic variation in coercion constructions across aspectual verbs. Results reveal that LMs struggle with retrieving plausible events and with ranking them above less plausible ones. Moreover, we found that LMs do not exploit the compositional properties of coercion sentences in their predictions.

1 Introduction

Sentences like *John began the book* are examples of complement coercion, a linguistic phenomenon displaying a type-mismatch between the required verb argument and the observed one (Pustejovsky, 1991, 1995): the aspectual verb (e.g., *begin*) semantically requires an event-denoting argument but is composed with an entity as its syntactic complement. Although the event is not overtly expressed, we are able to recover a plausible candidate exploiting contextual or lexical information (Pustejovsky, 1991, 1995; Lapata and Lascarides, 2003): the sentence above can then be interpreted as *John began (reading, writing) the book*.

Complement coercion has drawn attention as a potential violation of the Fregean principle of compositionality. Compositionality implies that all aspects of sentence meaning should originate in the meanings of the constituent parts, given the way those are combined syntactically (Asher, 2015). The interpretation of implicit elements results then from enriched composition (Jackendoff,

1997): semantic processes that exploit conceptual meaning, discourse context and world knowledge (Pustejovsky, 1991, 1995; Jackendoff, 1997).

Language Models (LM) based on the Transformer architecture (Vaswani et al., 2017) have demonstrated remarkable capabilities in a wide range of NLP tasks, including natural language understanding. Despite their success, few studies have focused on enriched composition phenomena such as complement coercion (Gu, 2022; Ye et al., 2022). Some studies have investigated LM performance framing complement coercion as an event retrieval task and demonstrating the challenges of recovering underlying semantic information from coercion sentences (Rambelli et al., 2020; Ye et al., 2022; Gietz and Beekhuizen, 2022; Gu, 2022; Im and Lee, 2024; Rambelli et al., 2024). However, most studies have been conducted in English, a language with low variability in the syntax of coercion constructions. As a consequence, little is known about the interplay of syntax and semantics in covert event retrieval in LMs: (how) do machines exploit compositional properties of coercion sentences to arrive at plausible interpretations?

The current study makes three contributions. First, to our knowledge, it is the first study of LMs on complement coercion that both uses a language other than English, in this case Norwegian, and that evaluates and compares a set of different LMs (autoencoders and autoregressive models). Second, we investigate the interaction between different aspectual verbs and post-verbal constituents in canonical syntactic constructions of coercion sentences. Norwegian shows some variation in how complement coercion is syntactically realized, and therefore allows us to probe whether LMs are sensitive to syntactic and compositional semantic properties of these constructions across aspectual verbs.

Finally, Norwegian is currently considered a low-resource language by various authors (Kummer-vold et al., 2022; Liu et al., 2023; Samuel et al.,

2024), and we are releasing our evaluation dataset for complement coercion resolution in Norwegian. Complement coercion with aspectual verbs is statistically rare in Norwegian corpora (see below): recovering implicit events could be challenging for a ‘data hungry’ technology such as LMs.

2 Related Work

2.1 Complement Coercion in Norwegian

Complement coercion has been studied in several high-resourced languages. Apart from English, we find studies on German (Rüd and Zarcone, 2011; Zarcone and Padó, 2011; Zarcone et al., 2012, 2014), French (Godard and Jayez, 1993; Pustejovsky and Bouillon, 1995), Dutch (Sweep, 2012), and Chinese (Hsu and Hsieh, 2013), while there has been little research on Norwegian. Spalek (2015) analyzed the verb of cessation *avslutte* (to conclude), comparing Norwegian with English, Spanish and German. Spalek concluded that coercion is limited to a reduced set of entities that can be combined with the verb, especially “information-content entities” (e.g, text) (Spalek, 2015, p. 531). Spalek and Sæbø (2019) argued that Norwegian speakers tend to combine eventive verbs with specific particles that denote a particular stage of the event (e.g., *å stryke ferdig*, to finish ironing).

Radaelli and Baggio (2025) conducted a study on the Norwegian Colossal Corpus (NCC) (Kummervold et al., 2022), a large set of corpora that includes approximately 21M documents for a total of 7B tokens. The study considered a wider class of aspectual verbs than previous theoretical research: *begynne* (to begin), *starte* (to start), *fortsette* (to continue), *ende* (to end), and *avslutte* (to conclude). The authors found that the syntax of complement coercion in Norwegian can vary according to the aspectual verb. Initiation verbs are usually combined with PPs introduced by the prepositions *på* or *med*:

- (1) *Gutten begynte/startet på/med boken.*
(The boy began/started [with] the book.)

These combinations appear with higher frequency in complement coercion sentences compared to other aspectual verbs. The continuation verb *fortsette* introduces coercion mainly with *med*-PPs, and, to a lesser extent, directly with nominals:

- (2) *Gutten fortsatte [med] boken.*
(The boy continued [with] the book.)

The cessation verb *avslutte* prefers direct objects, while *med*-prepositional phrases appear less often:

- (3) *Gutten avsluttet [med] boken.*
(The boy finished [with] the book.)

Not all aspectual verbs can trigger coercion phenomena (e.g., the verb *ende* was excluded), nor do aspectual verbs significantly differ in occurrence frequency in coercion constructions.

The corpus analysis confirms the findings of Spalek (2015) and Spalek and Sæbø (2019): complement coercion occurs with a restricted set of entity categories. Although a similar trend can be found also in other languages (e.g., see Verspoor (1997) for English and Rüd and Zarcone (2011) for German), Norwegian shows even less variability, reducing the set of entities primarily to everyday objects such as text, music, songs, food and drinks.

Considering Pustejovsky’s Generative Lexicon perspective (Pustejovsky, 1995), the productivity of coercion can also be limited by the interaction of syntactic and semantic factors. If, on the one hand, entities admit either AGENTIVE or TELIC qualia readings, their combination with prepositions may further reduce the set of plausible event candidates. The preposition *med* appears to play a ‘passe-partout’ role, with greater flexibility in event interpretation, admitting not only default qualia readings but also contextual information, if present. The preposition *på*, on the other hand, tends to further constrain interpretations: the corpus data showed a stronger tendency to express AGENTIVE interpretations with entities that are created rather than utilized. In a survey conducted in parallel to the corpus analysis, Radaelli & Baggio also found that Norwegian speakers prefer to express similar concepts to complement coercion through a broad range of phrasal constructions (e.g., *å sette i gang*, to begin). The study concluded that complement coercion is a relatively low-frequency phenomenon, with slightly more than 1500 cases over 79,000 sentences (approx. 1% of cases) with aspectual verbs and syntactic constructions compatible with coercion.

2.2 LM Approaches to Complement Coercion

In one of the first studies testing LMs on complement coercion, Rambelli et al. (2020) evaluated the events retrieved by pretrained models of the BERT and the GPT families. They found that LMs performed well, but not significantly better than the

best distributional models.

Ye et al. (2022) claimed that Transformer-based models can learn coercion interpretations via "dense paraphrasing" (DP): DP involves the reformulation of a given coercion sentence in a way that eventive information is revealed, ambiguity is removed and the original sentence meaning is preserved. They found that BERT struggles in interpreting coercion, but a fine-tuning with explicitly paraphrased sentences improved its performance.

Finally, Gu (2022) investigated the behavior of GPT-2 on complement coercion by analyzing surprisal estimates. The goal was to understand how LMs process coercion constructions at the VP. Significant surprisal effects were observed in the target region, aligning with psycholinguistic findings of increased processing costs at the complement (McElree et al., 2001; Traxler et al., 2002).

3 Experimental Settings

3.1 Task Proposal

In previous work on complement coercion, evaluations typically compared a discrete set of high likelihood predictions against a predefined set of gold standard outputs. In our study, we use a different evaluation approach, relying on a ranked prediction distribution rather than just the most probable outputs: for every context-neutral sentence¹ s belonging to a set S , a given model m generates a set of top-k-ranked output predictions $O = \{o_1 \dots o_k\}$. We then evaluate each output with the mean average precision metric, allowing us to determine to what extent LMs consistently predicts meaningful eventive interpretations in their ranking. The distribution should reflect a re-ranking of tokens when the model is exposed to coercion sentences, providing evidence of its sensitivity to coercion. In cases where a LM is exposed to a sentence such as *The boy began the book*, we expect that the combination of the triplet <subject, coercion verb, entity> would result in a re-ordering of implicit events (see Figure 1): the ranking should reflect the interaction of the triple's composition, where plausible verbs (events) are collocated at the top of the rank as the most likely interpretations.

Instead of using a set of predefined interpretations, our study will consider every eventive interpretation that meets the syntactic and semantic

constraints of complement coercion as correct. According to Piñango and Deo (2016) and Spalek and Sæbø (2019) the covert event of a complement coercion sentence should be telic: combined with the subject and complement, it must establish a natural endpoint or goal state.

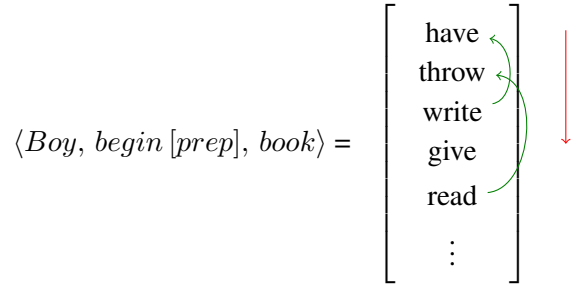


Figure 1: An example of re-ranking candidate events when the expressions in the given triplet are composed.

Given the above requirements, LMs predicted events should be evaluated considering their Aktionsart (lexical aspect) class by using Vendler's classification system (Vendler, 1967). We identify the class of **accomplishment** verbs as our ground truth in this task, as they denote dynamic and durative actions with a specific endpoint, aligning with the telicity criterion by Spalek and Sæbø (2019). As there is no predefined set of implicit events for interpreting a coercion sentence, we consider as compositionally plausible candidates all predicted verbs that fit within the accomplishment class. In case a model predicts events weakly associated with a specific coercion triplet (e.g., *begin the book* \rightarrow *eat*, see Lascarides and Copestake (1998)), this does not necessarily indicate low performance: the output can count as correct, if the retrieved event is an accomplishment. It is possible to construct contexts where even apparently deviant events are plausible, so long as they are accomplishments: e.g., *The goat began (eating) the book*.

3.2 Dataset

We created a new dataset with pairs with (a) a context-neutral sentence with a coercion triplet and varying syntactic structure (*på*-NP, *med*-NP, NP) and (b) a sentence prompting event resolution:

(a) Kim {VERB-FIN} {PREP|Ø} {ENTITY-DEF}.

(b) Det som Kim {VERB-FIN} å gjøre, var å [MASK]. (What Kim {VERB-FIN} to do,

¹Context-neutral, canonical coercion sentences contain the subject, the aspectual verb and its complement, with unmarked word order and no additional context.

was [MASK].)

Each placeholder in brackets is replaced with the relevant lexical item. The template encompasses a combination of the following elements:

- 90 entities ({ENTITY-DEF}) were carefully selected to represent real artifacts, avoiding abstract and ambiguous concepts. In addition, following Piñango and Deo (2016, p. 387), we ensured that entities can be semantically interpreted as "incremental theme arguments of the implicit event", a crucial element in coercion configurations. We included entities that never occurred in coercion sentences in the NCC corpus study, ensuring that the models (especially those trained exclusively on NCC) are exposed to sentences not seen during pre-training. Six distinct entity categories were used: *food*, *text*, *clothing*, *everyday objects*, *construction/housing*, and *entertainment*. Entities were only used in definite form.
- Four aspectual verbs ({VERB-FIN}), namely *begynne* (begin), *starte* (start), *fortsette* (continue), and *avslutte* (finish) were composed with each entity. The verb was always presented in the same simple past form (*preteritum*) in both sentences in a pair.
- Three syntactic constructions were used ({PREP|Ø}): the complement is either introduced by a PP with the prepositions **på** or **med** followed by the NP denoting an entity, or only by the latter NP.
- The same subject was used for every sentence, with a neutral name (*Kim*) to avoid gender and other biases that may affect the results.
- In all pairs, the prompt (b) included the [MASK] token the model has to predict.

A total of 1080 sentence pairs in standard written Bokmål form were used to each model.

3.3 Models

We evaluated a total of 17 different pre-trained Norwegian LMs varying in framework, parameter size, and training data. The models belong to two broad categories: BERT-like autoencoder models, and autoregressive models such as GPT-2 (Radford et al., 2019), LLAMA-2 (Touvron et al., 2023), Mistral

Model	# Par.	Tr. Data
mBERT CASED/UNCASED	178M	3.3B*
NB-BERT-BASE	178M	7B
NB-BERT-LARGE	355M	7B
NORBERT	111M	1.9B
NORBERT2	125M	15B
NORBERT3-base	123M	25B
NORBERT3-large	353M	25B
NORBERT3-SMALL	40M	25B
NORBERT3-XS	15M	25B
NORBLOOM-7B-SCRATCH	7B	26.7B
NORGPT-369M	369M	25B
NORGPT-3B	3B	25B
NORGPT-3B-CONTINUE	3B	25B
NORLLAMA-3B	3B	26.7B
NORMISTRAL-7B-SCRATCH	7B	26.7B
NORMISTRAL-7B-WARM	7B	26.7B

Table 1: Tested LMs with approximate information on number of parameters (#Par.) and training data (Tr. Data). *mBERT was trained on 114 languages.

(Jiang et al., 2023) and Bloom (Scao et al., 2023) (Table 1). All models are available on Huggingface².

3.4 Baseline Model

To assess event retrieval in complement coercion by LMs, it is necessary to find a baseline model, here provided by the NCC, an open-source corpus used for training most LMs in Norwegian³. For each entity in the dataset, we extracted the most probable verbs (events) associated with the entity. The extracted verbs were determined on the basis of the Pointwise Mutual Information (PMI) score, a metric evaluating the association strength between two words w_1 and w_2 (Church and Hanks, 1990):

$$\text{PMI}(w_1, w_2) = \log \frac{P(w_1, w_2)}{P(w_1) \cdot P(w_2)}$$

In our study, the score was calculated through the joint probability between each sentence predicate (event) and its object (entity) in the entire corpus.

3.5 Evaluation and Annotation

We used two common evaluation metrics. One is mean average precision (mAP) (see Manning et al. (2009, from p. 159) and Kotlerman et al. (2010)):

$$\text{mAP} = \frac{1}{S} \sum_{s=1}^S \text{AP}(s)$$

²<https://huggingface.co/>

³The National Library of Norway, the maintainer of the NCC corpus, has recently updated its distribution policies, limiting access to some subcorpora. This has caused a decrease in the corpus size from approximately 7B to 4.5B tokens. The present study was conducted on the previous corpus version. More information on <https://huggingface.co/datasets/NbAiLab/NCC>

Model	mAP	A1
NCC (Baseline)	0.59	0.47
NORGPT-369M	0.56	0.54
NORGPT-3B	0.48	0.42
NORGPT-3B-CONTINUE	0.46	0.42
NORLLAMA-3B	0.71	0.67
BERT-BASE-MULTILINGUAL-CASED	0.07	0.00
BERT-BASE-MULTILINGUAL-UNCASED	0.27	0.22
NB-BERT-BASE	0.38	0.33
NB-BERT-LARGE	0.54	0.47
NORBERT	0.25	0.18
NORBERT2	0.44	0.34
NORBERT3-BASE	0.63	0.58
NORBERT3-LARGE	0.60	0.55
NORBERT3-SMALL	0.59	0.55
NORBERT3-XS	0.29	0.16
NORBLOOM-7B-SCRATCH	0.46	0.34
NORMISTRAL-7B-SCRATCH	0.38	0.29
NORMISTRAL-7B-WARM	0.63	0.54

Table 2: Mean average precision (mAP) and top-rank accuracy (A1) results on the covert event retrieval task in Norwegian.

It consists of the weighted means of average precision (AP) scores across all sentences (S):

$$AP(q) = \sum_{k=1}^5 P(k) \cdot \Delta R(k)$$

where P is the precision value calculated at the cut-off rank k and $\Delta R(k)$ is the change in recall (R) from rank $k - 1$ to k . mAP provides the ranking direction of models when complement coercion occurs. A high mAP value indicates a model that mostly considers accomplishment verbs in the prediction list, collocating them at the top, whereas a low mAP value suggests a failure in prioritizing accomplishment verbs as completions.

The second metric is the mean top-ranked accuracy (A1) in all sentences, considering only the most likely prediction in the ranking. This metric allows us to study what types of verb (events) the models consider as the most salient ones.

4 Results and Task Discussion

Table 2 shows the performance results of all LMs on the covert event retrieval task in Norwegian, with mAP and A1 scores. Model performance varies according to the interplay of two main factors: model framework and model size (number of parameters and training data). The NORBERT3 family shows relatively high performance compared to other BERT-like frameworks, with NORBERT3-BASE and NORBERT3-LARGE outperforming the baseline on both measures. Larger LMs outperform NB-BERT models and the previous generations of

NORBERT models, which showed poorer performances, possibly due to less training data available. Models like NORBERT3-XS performed less well probably due to their reduced parameter size despite the same amount of training data.

Almost all GPT-2-based models, as well as NORBLOOM-7B-SCRATCH and NORMISTRAL-7B-SCRATCH performed poorly, ranking below the baseline, despite their size. NORMISTRAL-7B-WARM outperforms the baseline in both cases, compared to the version trained from scratch: pretraining on the English vanilla version and successive pretraining on Norwegian data may have given the model an advantage, allowing for the transfer of rich representations from English text. Finally, NORLLAMA-3B can be considered as the most capable model among those tested here. Its success could be attributed perhaps to its large training corpus, with more than 25B training tokens in Norwegian and other Scandinavian languages.

Language models generally struggle to interpret complement coercion sentences. Overall low mAP scores suggest difficulties in generating plausible accomplishments among high-ranked candidate mask replacements. This is confirmed when cross-analyzing A1 scores: even the best model, NORLLAMA-3B fails to reach a 70% level of accuracy, indicating that a consistent chance that it non-accomplishments and other implausible verbs will be predicted as plausible interpretations at the top of the list. Similarly, the top-10 ranked models achieve an A1 score ranging from 0.42 to 0.58, indicating that they have on average 50% chance of failing to rank accomplishments at the top.

We will now turn to an analysis of model performance taking into account both mAP and A1 scores subdivided according to aspectual verbs and their syntactic structures in coercion sentences. For the sake of simplicity, we will consider the best performing model NORLLAMA-3B. The results are shown in Table 3.

Consistently high mAP scores are found with initiation verbs. The verb *starte* shows high mAP scores reaching 0.81 precision when entity arguments in coercion sentences are introduced by the preposition *med*, 0.79 with nominals, and 0.76 with the preposition *på*. The verb *begynne* was associated with worse performance, while showing a similar trend as *starte*. Sentences with entity arguments introduced by *med* reached 0.79 precision, 0.75 with nominals, and 0.73 with *på*. The two remaining aspectual verbs showed similar results,

NorLlama-3B verb	prep	mAP	A1
avslutte	∅	0.66	0.61
	med	0.75	0.69
	på	0.64	0.53
begynne	∅	0.75	0.72
	med	0.79	0.81
	på	0.73	0.71
fortsette	∅	0.64	0.57
	med	0.64	0.56
	på	0.59	0.46
starte	∅	0.79	0.80
	med	0.81	0.83
	på	0.76	0.76

Table 3: Mean average precision (mAP) and top-rank accuracy (A1) results of NORLLAMA-3B categorized by aspectual verbs (*begynne*, *starte*, *fortsette*, and *avslutte*) and syntactic composition (introduced either by prepositions *på* or *med*, or by nominal, here \emptyset) in coercion sentences.

and arguments with *med* as preposition obtained higher precision scores. In sentences with *fortsette*, both nominals and *med*-prepositional phrases reached the same score (0.64).

A1 scores show a similar trend. The model performs better when coercion sentences are introduced by *starte*, with 0.83 of A1 accuracy when the entity NP is introduced by *med*-PP, 0.80 without a preposition, and 0.76 by *på*-PP. The verb *begynne* also serves as trigger for complement coercion, with a A1 score of 0.81 with *med*-prepositional phrases, 0.72 with simple nominals, and 0.71 with verb argument phrases introduced by *på*.

Two key observations are suggested by this analysis. First, different aspectual verbs are associated with differences in model performance. Our results indicate that the model can recover the implicit meaning more easily with initiation verbs in coercion sentences. This is consistent with the corpus analysis of Radaelli and Baggio (2025), which showed that among all aspectual verbs, initiation verbs feature more frequently in coercion sentences. Second, we only find weak differences in performance as a function of the syntax of post-verbal arguments. This suggests that the type of syntactic structures in complement coercion sentences plays only a minor role in the model’s process of recovery of implicit meaning.

4.1 Model surprisal

Previous studies (see above) indicated that LMs struggle to consistently retrieve covert events for complement coercion sentences. To understand the reasons behind these prediction difficulties, we can study the model’s behavior when it is exposed to complement coercion sentences. We conducted a further analysis that complements the previous ranking results by computing surprisal estimates for coercion sentences. Surprisal is used in NLP and psycholinguistic studies to quantify effort during sentence processing (see Hale (2001); Levy (2008); Smith and Levy (2013); Salicchi et al. (2023); Oh and Schuler (2023); Shain et al. (2024)):

$$S(w_i) = -\log_2 P(w_i | w_1, \dots, w_{i-1})$$

Surprisal measures how unexpected a given word (w_i) is, given its left context ($w_1 \dots w_{i-1}$). Higher surprisal values indicate greater processing difficulty, as upcoming words are less predictable.

We will use surprisal to determine whether models expect an entity-denoting noun in coercion triples. Specifically, we will compare surprisal estimates for complement coercion sentences (e.g., *Kim begynte på boken*, Kim began (on) the book) with their overt eventive counterparts (e.g., *Kim leste boken*, Kim read the book). The events were selected considering the highest PMI scores between each accomplishment and its associated entity. In total, we examined 2,160 sentences, using the same sentences from the previous task (1,080 coercion, 1,080 overt) combining all aspectual verbs, all entities, and the same three different syntactic structures. To compute surprisal estimates, we used log-probabilities provided by model logits. As coercion and overt sentences may differ in length, we will compute surprisal estimates for sentences as the mean of each word’s surprisal:

$$S_{mean}(s) = \frac{1}{N} \sum_{i=1}^N S(w_i)$$

where N corresponds to the number of token present in a sentence s . Here too, we tested surprisal for NORLLAMA-3B as a high performance model in this task. For the calculation of surprisal estimates, we used the tool *minicons* on Python. The data were analyzed using Wilcoxon signed-rank tests to compare the surprisal values between coercion and overt sentences. We hypothesized that the model would show higher surprisal values

	Baseline model Coefficient (β)	Model With Syntax Coefficient (β)
Intercept (Nominals)	9.7983 ($p < 0.001$)	2.0463 ($p < 0.001$)
Coercion	0.9429 ($p < 0.001$)	1.4946 ($p < 0.001$)
Sequence Length	-1.0612 ($p < 0.001$)	1.3388 ($p < 0.001$)
Explicit	—	0.5517 ($p < 0.001$)
Med	—	-2.5418 ($p < 0.001$)
På	—	-2.2583 ($p < 0.001$)

Table 4: Effect of Syntax on Surprisal. Comparison of two models.

for coercion sentences than for overt ones. The results confirmed the hypothesis, showing a statistical difference in surprisal ($W = 367176$, $p < 0.001$). This suggests a tendency of the model to assign prediction logits with lower probabilities when processing coercion sentences.

In order to analyze the extent to which syntactic structure can influence surprisal in coercion sentences, we proposed and compared two regression models. As baseline, we ran a model on surprisal using only sentence type (coercion vs explicit) and sequence length as predictors. The second model also included syntactic structure as a predictor (with *på*-PPs, *med*-PPs and direct nominals as levels). The baseline model ($R^2 = 0.173$) revealed that coercion sentences significantly increased surprisal. Moreover, sequence length negatively correlated with surprisal, meaning that longer sentences led to lower surprisal values. The second model ($R^2 = 0.181$) shows a significant positive trend in the coercion condition, as the baseline model. On the other hand, sequence length shows in this case a positive effect on surprisal. Sentences with *med*-prepositional phrases demonstrate lowest surprisal, while sentences with *på* exhibit slightly higher surprisal, but still lower than in the nominal conditions. Comparing the variance of the two models ($\Delta R^2 = +0.008$), we find small improvements attributable to syntax. Prepositions therefore reduce surprisal in contrast to sentences with direct nominals, where *med*-sentences led to lower surprisal, followed by the *på*. A boxplot is shown in Figure 2, displaying the surprisal values for the different syntactic structures.

5 General Discussion and Conclusion

The analyses carried out in this study clearly show that complement coercion remains an open challenge for LMs in low-resource languages, such as Norwegian. We investigated the extent to which LMs could recover implicit events in complement

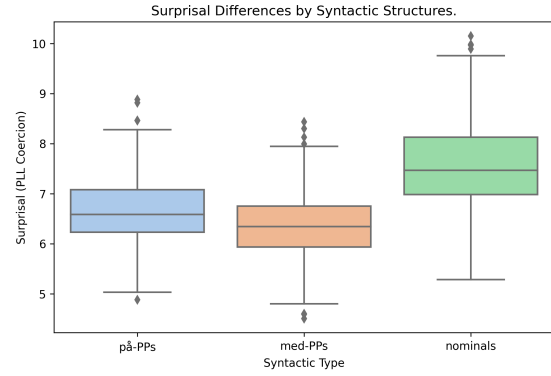


Figure 2: Boxplot describing surprisal values across syntactic structures.

coercion sentences. If models recognize these as coercion constructions, that require event retrieval, they should be able to distribute verb (event) predictions in such a way that accomplishments are ranked as the most probable covert events.

However, the outcomes of the event retrieval task indicate that LMs still have difficulties recovering implicit events. In particular, A1 scores are consistently low across models, suggesting a failure to retrieve potential accomplishment verbs as the most likely event predictions in the task. Moreover, the mAP scores confirmed the models’ limitations, as they fail to systematically rank accomplishment verbs in higher positions consistently. Only few models could outperform the baseline, whose predictions are based on simple statistical calculations on the NCC corpus frequency: this is significant, considering that such models were trained on corpora 3.5 times bigger than the baseline size. The results also highlight performance differences across model architectures:

- NORLLAMA-3B outperformed all the models, obtaining a moderate performance score; its success may be due to its new improved architecture and training optimizations (e.g., SwiGLU activation function, Grouped Query Attention mechanism, rotary positional embeddings) combined with a high amount of training data.
- GPT-models, on the other hand, could not perform the task efficiently, even with the larger models. Probably, the traditional autoregressive GPT-2 lacks an efficient architecture to capture underlying information like covert events evoked by coercion phenomena.
- Even NORBLOOM-7B-SCRATCH and

563
564
565
566
567
568
569
570
571
572
573
574
575
576

577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599

600
601
602
603
604
605
606
607
608
609
610
611
612
613

NORMISTRAL-7B-SCRATCH performed poorly despite their size. Their low performance could be attributed to its training phase done solely on Norwegian, especially when compared with the best performing NORMISTRAL-7B-WARM with a prior pre-training phase on English. Moreover, NORBLOOM-7B-SCRATCH derives from a 176B-parameter multilingual model (Scao et al., 2023). We can suppose that a model rescaling to 7B parameters pretraining solely on Norwegian could have led to undertraining and to failure in capturing semantic and syntactic information in an effective way.

- The NORBERT3 family, especially the base and large versions, could reach moderate performance levels despite their reduced number of parameters. The BERT architecture appears to be well-suited for learning and storing world knowledge and relational knowledge between words during pretraining, making them effective for cloze tasks (Petroni et al., 2019; Rogers et al., 2020). In addition, their customized autoencoder framework, incorporating the extended MLM pre-training task (Samuel et al., 2023), may have facilitated acquisition of syntactic and semantic information relevant for the present task. To this purpose, it should be noticed that also in the complement coercion study of Rambelli et al. (2020) on English a bidirectional architecture (RoBERTa) was the one showing the highest correlations with human production frequencies for the candidate covert event. However, vanilla architectures combined with less training data would drastically reduce performance as seen in the NB-BERT models.

To better understand how models process coercion sentences and investigate the reasons behind their difficulties in event retrieval task, we compared surprisal estimates between coercion sentences and their overt event counterparts. Higher surprisal values for complement coercion sentences suggest that LMs generally find coercion constructions less predictable, which should be expected given their relative infrequency in Norwegian corpora. However, rare constructions in human language can still be interpreted by exploiting lexical meaning and syntactic structure, even when context is minimal or absent: overall, our results suggest that LMs are unable to exploit available composi-

tional information to generate accomplishments as plausible event completions in complement coercion sentences.

A more detailed analysis of the best performing model revealed only moderate variation in performance according to the specific aspectual verb used. Initiation verbs lead to better performance: based on results of corpus studies, this is probably due to stronger statistical associations between these aspectual verbs and (particular classes of) entity-denoting nominals. However, we could not find clear differences between different syntactic constructions within the same aspectual verbs, which suggests that models cannot exploit differences in syntactic structure to recognize these as coercion constructions and accordingly attempt the retrieval of plausible accomplishments.

Linear regression models were also proposed for checking whether coercion surprisal estimates were influenced by the syntactic structures proposed in the dataset. Results revealed weak differences in surprisal estimates, especially between coercion sentences with entity-denoting complements introduced by prepositions or directly by NPs, showing greater processing difficulties in the latter cases. This partially aligns with the results presented in table 3, where nominals led to lower scores, while *med*-PPs were associated with better performance. Moreover, LM behavior aligns weakly with the NCC corpus study by Radaelli and Baggio (2025): the authors found that *med*-prepositional phrases occur more frequently in coercion constructions and allow greater flexibility in event interpretations.

Considering LM’s inability to exploit compositionality (lexical meaning and syntactic structure) with complement coercion sentences, future work should explore what other factors can impact LM’s performance in this task. There are at least two possible research directions. First, an analysis of the role of linguistic context as a factor in performance improvement: what aspects of sentence or discourse context can facilitate event retrieval? Second, an analysis of the extent to which LM’s performance is dependent on ontology of entity types: can event retrieval be facilitated by specific classes of entities, as suggested by theoretical linguistic and corpus research?

References

Nicholas Asher. 2015. *Types, meanings and coercions in lexical semantics*. *Lingua*, 157:66–82. Polysemy:

614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660

661
662
663

664	Current Perspectives and Approaches.	
665	Kenneth Church and Patrick Hanks. 1990. Word associ-	
666	ation norms, mutual information, and lexicography.	
667	<i>Computational Linguistics</i> , 16(1):22–29.	
668	Frederick G Gietz and Barend Beekhuizen. 2022. Re-	
669	modelling complement coercion interpretation. <i>Soci-</i>	
670	<i>ety for Computation in Linguistics</i> , 5(1).	
671	Daniele Godard and Jacques Jayez. 1993. Towards a	
672	proper treatment of coercion phenomena . In <i>Sixth</i>	
673	<i>Conference of the European Chapter of the Associ-</i>	
674	<i>ation for Computational Linguistics</i> , Utrecht, The	
675	Netherlands. Association for Computational Linguis-	
676	tics.	
677	Yuling Gu. 2022. Measure more, question more: Ex-	
678	perimental studies on transformer-based language	
679	models and complement coercion. <i>arXiv preprint</i>	
680	<i>arXiv:2212.10536</i> .	
681	John Hale. 2001. A probabilistic earley parser as a	
682	psycholinguistic model. In <i>Proceedings of NAACL</i> .	
683	Chan-Chia Hsu and Shu-Kai Hsieh. 2013. To Coerce or	
684	Not to Coerce: A Corpus-based Exploration of Some	
685	Complement Coercion Verbs in Chinese . In <i>Proceed-</i>	
686	<i>ings of the 6th International Conference on Gener-</i>	
687	<i>ative Approaches to the Lexicon (GL2013)</i> , pages	
688	13–20, Pisa, Italy. Association for Computational	
689	Linguistics.	
690	Seohyun Im and Chungmin Lee. 2024. What gpt-4	
691	knows about aspectual coercion: Focused on “begin	
692	the book”. In <i>Proceedings of the Workshop on Cogni-</i>	
693	<i>tive Aspects of the Lexicon@ LREC-COLING 2024</i> ,	
694	pages 56–67.	
695	Ray Jackendoff. 1997. <i>The Architecture of the Lan-</i>	
696	<i>guage Faculty</i> . MIT Press. Google-Books-ID:	
697	_AkHjNi6qqcC .	
698	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	
699	sch, Chris Bamford, Devendra Singh Chaplot, Diego	
700	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	
701	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	
702	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	
703	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	
704	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,	
705	<i>arXiv:2310.06825</i> .	
706	Lili Kotlerman, Ido Dagan, Idan Szpektor, and Maayan	
707	Zhitomirsky-Geffet. 2010. Directional distributional	
708	similarity for lexical inference. <i>Natural Language</i>	
709	<i>Engineering</i> , 16(4):359–389.	
710	Per Kummervold, Freddy Wetjen, and Javier de la Rosa.	
711	2022. The Norwegian Colossal Corpus: A Text Cor-	
712	pus for Training Large Norwegian Language Models .	
713	In <i>Proceedings of the Thirteenth Language Resources</i>	
714	<i>and Evaluation Conference</i> , pages 3852–3860, Mar-	
715	seille, France. European Language Resources Asso-	
716	ciation.	
	Maria Lapata and Alex Lascarides. 2003. A Probabilis-	717
	tic Account of Logical Metonymy . <i>Computational</i>	718
	<i>Linguistics</i> , 29(2):261–315.	719
	Alex Lascarides and Ann Copestake. 1998. Pragmatics	720
	and word meaning. <i>Journal of linguistics</i> , 34(2):387–	721
	414.	722
	Roger Levy. 2008. Expectation-based syntactic compre-	723
	hension. <i>Cognition</i> , 106(3):1126–1177.	724
	Peng Liu, Lemei Zhang, Terje Farup, Even W Lauvrak,	725
	Jon Espen Ingvaldsen, Simen Eide, Jon Atle Gulla,	726
	and Zhirong Yang. 2023. Nlebench+ norglm: A	727
	comprehensive empirical analysis and benchmark	728
	dataset for generative language models in norwegian.	729
	<i>arXiv preprint arXiv:2312.01314</i> .	730
	Christopher D. Manning, Prabhakar Raghavan, and Hin-	731
	rich Sch��tze. 2009. <i>An Introduction to Information</i>	732
	<i>Retrieval</i> . Cambridge University Press, Cambridge,	733
	England.	734
	Brian McElree, Matthew J Traxler, Martin J Pickering,	735
	Rachel E Seely, and Ray Jackendoff. 2001. Reading	736
	time evidence for enriched composition. <i>Cognition</i> ,	737
	78(1):B17–B25.	738
	Byung-Doh Oh and William Schuler. 2023.	739
	Transformer-based language model surprisal	740
	predicts human reading times best with about two	741
	billion training tokens. In <i>Findings of EMNLP</i> .	742
	Fabio Petroni, Tim Rockt��schel, Patrick Lewis, An-	743
	ton Bakhtin, Yuxiang Wu, Alexander H Miller, and	744
	Sebastian Riedel. 2019. Language models as knowl-	745
	edge bases? <i>arXiv preprint arXiv:1909.01066</i> .	746
	Maria Mercedes Pi��ango and Ashwini Deo. 2016. Re-	747
	analyzing the Complement Coercion Effect through	748
	a Generalized Lexical Semantics for Aspectual Verbs .	749
	<i>Journal of Semantics</i> , 33(2):359–408.	750
	James Pustejovsky. 1991. The Generative Lexicon .	751
	<i>Computational Linguistics</i> , 17(4):409–441.	752
	James Pustejovsky. 1995. <i>The Generative Lexicon</i> . MIT	753
	Press. Google-Books-ID: p2GJBnpBVg4C.	754
	James Pustejovsky and Pierrette Bouillon. 1995. As-	755
	pectual Coercion and Logical Polysemy . <i>Journal of</i>	756
	<i>Semantics</i> , 12(2):133–162.	757
	Matteo Radaelli and Giosu�� Baggio. 2025. Comple-	758
	ment Coercion with Aspectual Verbs Is Statistically	759
	Infrequent in Written Norwegian (forthcoming).	760
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	761
	Dario Amodei, and Sutskever. 2019. Language Mod-	762
	els are Unsupervised Multitask Learners . <i>OpenAI</i>	763
	<i>blog</i> , 1(8):9.	764
	Giulia Rambelli, Emmanuele Chersoni, Alessandro	765
	Lenci, Philippe Blache, Chu-Ren Huang, and 1	766
	others. 2020. Comparing Probabilistic, Distribu-	767
	tional and Transformer-based Models on Logical	768
	Metonymy Interpretation. In <i>Proceedings of AACL-</i>	769
	<i>IJCNLP</i> .	770

771	Giulia Rambelli, Emmanuele Chersoni, Davide Testa,	Josefien Sweep. 2012. Logical Metonymy in Dutch and	825
772	Philippe Blache, and Alessandro Lenci. 2024. Neural	German: Equivalents of Begin, Finish, and Enjoy.	826
773	Generative Models and the Parallel Architecture of	<i>International Journal of Lexicography</i> , 25(2):117–	827
774	Language: A Critical Review and Outlook. <i>Topics in</i>	151.	828
775	<i>Cognitive Science</i> .		
776	Anna Rogers, Olga Kovaleva, and Anna Rumshisky.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	829
777	2020. A primer in bertology: What we know about	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	830
778	how BERT works . <i>CoRR</i> , abs/2002.12327.	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	831
779		Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	832
780	Stefan Rüd and Alessandra Zarcone. 2011. Covert	Grave, and Guillaume Lample. 2023. LLaMA: Open	833
781	events and qualia structures for german verbs. In	and Efficient Foundation Language Models . <i>arXiv</i>	834
782	<i>Proceedings of the Metonymy 2011 Workshop</i> , pages	<i>preprint</i> . ArXiv:2302.13971 [cs].	835
783	17–22.		
784	Lavinia Salicchi, Emmanuele Chersoni, and Alessandro	Matthew J Traxler, Martin J Pickering, and Brian McEl-	836
785	Lenci. 2023. A study on surprisal and semantic relat-	ree. 2002. Coercion in sentence processing: Evi-	837
786	edness for eye-tracking data prediction. <i>Frontiers in</i>	dence from eye-movements and self-paced reading.	838
787	<i>Psychology</i> , 14.	<i>Journal of Memory and Language</i> , 47(4):530–547.	839
788			
789	David Samuel, Andrey Kutuzov, Samia Touileb, Erik	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	840
790	Velldal, Lilja Øvrelid, Egil Rønningstad, Elina Sigdel,	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	841
791	and Anna Palatkina. 2023. NorBench – A Bench-	Kaiser, and Illia Polosukhin. 2017. Attention is all	842
792	mark for Norwegian Language Models . In <i>Proceed-</i>	you need. <i>Advances in neural information processing</i>	843
793	<i>ings of the 24th Nordic Conference on Computational</i>	<i>systems</i> , 30.	844
794	<i>Linguistics (NoDaLiDa)</i> , pages 618–633, Tórshavn,		
795	Faroe Islands. University of Tartu Library.	Z Vendler. 1967. <i>Linguistics in philosophy</i> ithaca, ny:	845
796		Cornell univ.	846
797	David Samuel, Vladislav Mikhailov, Erik Velldal, Lilja		
798	Øvrelid, Lucas Georges Gabriel Charpentier, and	Cornelia Maria Verspoor. 1997. Conventionality-	847
799	Andrey Kutuzov. 2024. Small languages, big models:	governed logical metonymy . In <i>Proceedings of the</i>	848
800	A study of continual training on languages of norway.	<i>second international workshop on computational se-</i>	849
801	<i>arXiv preprint arXiv:2412.06484</i> .	<i>mantics</i> , pages 300–312. Citeseer.	850
802			
803	Teven Le Scao, Angela Fan, Christopher Akiki, El-	Bingyang Ye, Jingxuan Tu, Elisabetta Jezek, and James	851
804	lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman	Pustejovsky. 2022. Interpreting logical metonymy	852
805	Castagné, Alexandra Sasha Luccioni, François Yvon,	through dense paraphrasing. In <i>Proceedings of the</i>	853
806	Matthias Gallé, Jonathan Tow, Alexander M. Rush,	<i>Annual Meeting of the Cognitive Science Society</i> ,	854
807	Stella Biderman, Albert Webson, Pawan Sasanka Am-	volume 44.	855
808	manamanchi, Thomas Wang, Benoît Sagot, Niklas		
809	Muennighoff, Albert Villanova del Moral, and 373	Alessandra Zarcone and Sebastian Padó. 2011. Gener-	856
810	others. 2023. Bloom: A 176b-parameter open-	alized Event Knowledge in Logical Metonymy Reso-	857
811	access multilingual language model . <i>Preprint</i> ,	lution.	858
812	arXiv:2211.05100.		
813		Alessandra Zarcone, Sebastian Padó, and Alessandro	859
814	Cory Shain, Clara Meister, Tiago Pimentel, Ryan Cot-	Lenci. 2014. Logical Metonymy Resolution	860
815	terrell, and Roger Levy. 2024. Large-scale evidence	in a Words-as-Cues Framework: Evidence	861
816	for logarithmic effects of word predictability on read-	From Self-Paced Reading and Probe Recogni-	862
817	ing time. <i>Proceedings of the National Academy of</i>	<i>tion</i> . <i>Cognitive Science</i> , 38(5):973–996. <i>_eprint</i> :	863
818	<i>Sciences</i> , 121(10):e2307876121.	https://onlinelibrary.wiley.com/doi/pdf/10.1111/cogs.12108 .	864
819			
820	Nathaniel J Smith and Roger Levy. 2013. The effect	Alessandra Zarcone, Jason Utt, and Sebastian Padó.	865
821	of word predictability on reading time is logarithmic.	2012. Modeling covert event retrieval in logical	866
822	<i>Cognition</i> , 128(3):302–319.	metonymy: probabilistic and distributional accounts.	867
823		In <i>Proceedings of the 3rd workshop on cognitive mod-</i>	868
824	Alexandra Anna Spalek. 2015. The Influence of Con-	<i>eling and computational linguistics (CMCL 2012)</i> ,	869
825	text in Meaning: The Panorama of Complement Co-	pages 70–79.	870
826	ercion . In <i>Modeling and Using Context</i> , pages 526–		
827	531, Cham. Springer International Publishing.		
828			
829	Alexandra Anna Spalek and Kjell Johan Sæbø. 2019.		
830	To Finish in German and Mainland Scandinavian:		
831	Telicity and Incrementality . <i>Journal of Semantics</i> ,		
832	36(2):349–375.		