

# Sequential Compositional Generalization in Multimodal Models

Semih Yagcioglu<sup>1</sup> Osman Batur İnce<sup>2,3</sup> Aykut Erdem<sup>2,3</sup>

Erkut Erdem<sup>1</sup> Desmond Elliott<sup>4,5</sup> Deniz Yuret<sup>2,3</sup>

<sup>1</sup>Hacettepe University <sup>2</sup>Koç University <sup>3</sup>KUIS AI Center

<sup>4</sup>University of Copenhagen <sup>5</sup>Pioneer Centre for AI

## Abstract

The rise of large-scale multimodal models has paved the pathway for groundbreaking advances in generative modeling and reasoning, unlocking transformative applications in a variety of complex tasks. However, a pressing question that remains is their genuine capability for stronger forms of generalization, which has been largely underexplored in the multimodal setting. Our study aims to address this by examining sequential compositional generalization using COMPACT (Compositional Activities)<sup>1</sup>, a carefully constructed, perceptually grounded dataset set within a rich backdrop of egocentric kitchen activity videos. Each instance in our dataset is represented with a combination of raw video footage, naturally occurring sound, and crowd-sourced step-by-step descriptions. More importantly, our setup ensures that the individual concepts are consistently distributed across training and evaluation sets, while their compositions are novel in the evaluation set. We conduct a comprehensive assessment of several unimodal and multimodal models. Our findings reveal that bi-modal and tri-modal models exhibit a clear edge over their text-only counterparts. This highlights the importance of multimodality while charting a trajectory for future research in this domain.

## 1 Introduction

Humans possess a remarkable ability to rapidly understand new concepts by leveraging and combining prior knowledge. This compositional generalization allows for an understanding of complex inputs as a function of their constituent parts. For instance, having grasped the meanings of “dax” and “walk twice” humans can effortlessly understand “dax twice” (Lake and Baroni, 2018). However, even as neural networks trained on increasingly larger datasets achieve impressive results across a wide range of tasks, their ability to compositionally

generalize remains limited. Recently, the research community has demonstrated growing interest in evaluating models under different distributions, such as temporal shifts (Lazaridou et al., 2021; Liska et al., 2022), or unseen compositions (Lake and Baroni, 2018; Ettinger et al., 2018; Bahdanau et al., 2019; Surís et al., 2020). Within the domain of multimodal learning, prior investigations into compositionality have primarily delved into visual grounding (Thrush et al., 2022), downstream multimodal tasks like image captioning (Nikolaus et al., 2019; Jin et al., 2020) and visual question answering (Bahdanau et al., 2019), or vocabulary acquisition from videos (Surís et al., 2020) or with interactive agents (Hill et al., 2019).

Addressing the challenge of compositional generalization in the context of multimodal models is increasingly important with the recent advances in large multimodal foundation models, such as GPT-4 (OpenAI, 2023), Flamingo (Alayrac et al., 2022), and IDEFICS (Laurençon et al., 2023). Experimenting with closed-source or proprietary models introduces challenges, including reproducibility issues, associated costs, and limited transparency regarding their development and training methodologies (Nityasya et al., 2023). This inspires us to investigate the potential of open-source large multimodal foundation models – IDEFICS in particular, for multimodal sequential compositional generalization, which we define as the model’s capability to understand and generate predictions about novel compositions of primitive elements derived from sequential multimodal inputs – for instance, video data wherein actions unfold in a discernible order. Consider the process of cooking onions: one typically needs to *peel* and *slice* an ONION before *frying* it in a PAN. Our central inquiry revolves around the proficiency of models in comprehending such sequential and compositional activities.<sup>2</sup>

<sup>1</sup>Project Page: <http://cyberiada.github.io/CompAct>

<sup>2</sup>Note that this differs from in-context learning, where


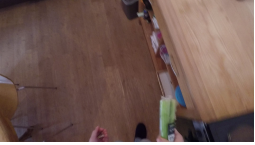
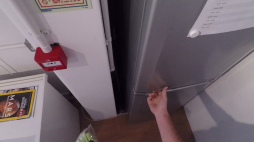







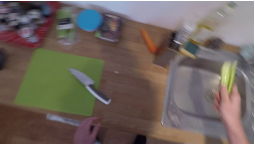
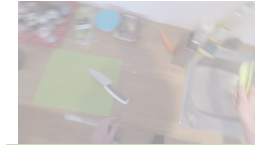
	Inputs (keyframes and utterances)		Targets (next utterance)	
Training		take celery		throw things into garbage bin
		open fridge		put celery back into fridge
		pour sesame oil		close sesame oil
		pick up onion		cut onion in half
Evaluation		wash celery		close tap
		put down celery		cut celery

Figure 1: Overview of the compositional generalization setup in our COMPACT dataset. During training, the model has seen the verbs *wash*, *close*, *put down*, *throw*, *open*, *pour*, *cut* and *pick up* with the objects GARBAGE BIN, FRIDGE, SESAME OIL, ONION, and CELERY. It has never seen the composition of *cut* and CELERY, and thus needs to generalize to this novel composition at test time.

In this study, we introduce COMPACT (Compositional Activities) to investigate multi-modal sequential compositional generalization, a uniquely constructed compositional dataset curated from the EPIC KITCHENS-100 dataset (Damen et al., 2022, EK-100). The EK-100 dataset encompasses 100 hours of egocentric video footage from 45 distinct kitchens, documenting people performing routine household tasks. Each video contains three streams of information: *visual data* in the videos; *audio data* involving non-narrative audio elements –such as the sounds associated with chopping an onion; and *textual data* in the form of short, crowd-sourced descriptions of the depicted activities, like “slice the carrot”, “pick up the milk”, or “wash the plate”. From these descriptions, individual verb and object concepts such as *slice*, *pick up*, *wash*, and CARROT, MILK, PLATE are extracted. The compositional splits are devised based on the verb and object concepts gleaned from the video descriptions, resulting in training and evaluation sets showcasing similar distributions of atomic concepts but featuring varied combinations therein. Consequently, models should compositionally generalize from the training data. Aligning with the “dax twice” principle from Lake and Baroni (2018), if a model has been

trained with videos illustrating how to *slice* various food items, excluding ROOT VEGETABLES, then it should be capable of compositionally generalizing to understand what it means to *slice* the ROOT VEGETABLES from previously unseen instances.

In our study, we conduct a comprehensive evaluation of publicly available models, encompassing encoder-only pretrained models such as ImageBind (Girdhar et al., 2023) and MERLOT Reserve (Zellers et al., 2022) in addition to (multimodal) large language models (LLMs) like LLaMA2 (Touvron et al., 2023) and IDEFICS (Laurençon et al., 2023). These models exhibit versatility in processing various combinations of input streams, ranging from language-only to combinations like video + language, video + audio, and even video + language + audio. Our key experimental finding indicates the formidable challenge that all of these models face in mastering compositional generalization. Yet, it becomes abundantly clear that the utilization of multimodal input sources yields discernible advantages, suggesting a promising direction for refining future models.

## 2 The COMPACT Dataset

In our pursuit to systematically examine multi-modal sequential compositional generalization, we devised the COMPACT dataset, leveraging sequences from the EK-100 dataset (Damen et al.,

large-scale pretrained models are prompted for a task in a zero-shot setting, given a support set of task demonstrations.

2022). As previously noted, each video in the EK-100 features first-person perspectives of unscripted kitchen activities occurring within natural household environments. A video is composed of a sequence of shorter clips, represented as  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_k)$ , each of which is accompanied by manually annotated English narrations, denoted by  $\mathbf{x}_1, \dots, \mathbf{x}_k$ , describing the activities within. Additionally, these clips are integrated with audio tracks,  $\mathbf{a}_1, \dots, \mathbf{a}_k$ , which contain the sounds of occurring actions. We define an *instance* in the dataset as a combination of video – audio – narration:  $(\mathbf{V}, \mathbf{A}, \mathbf{X})$ . Each instance consists of a window of 4 clips, with the initial 3 clips serving as context and the last one designated for prediction.

Given this dataset, our primary focus is to facilitate researchers in exploring how multimodal models compositionally generalize to unseen combinations of concepts. We meticulously curate the COMPACT dataset to ensure a specific property: the individual concepts are consistently distributed across training and evaluation sets, while their compositions are novel in the evaluation set. This design mandates that a model should exhibit systematic generalization when interpreting the evaluation set. To illustrate, refer to the example shown in Fig. 1. During training, the model comes across nouns such as CELERY, GARBAGE BIN, FRIDGE, ONION, and verbs including *take*, *wash*, *close*, *put down*. In our evaluation set, we seek instances where an object-verb composition has not been previously encountered during training; for example, the unique pairing of the *cut* with the CELERY.

## 2.1 Forming the Compositional Splits

We use the *Maximum Compound Divergence* heuristic (Keysers et al., 2020) to curate a dataset that requires compositional generalization. The EK-100 dataset is annotated with 97 verb classes and 300 noun classes; these become the noun and verb *atoms*. Each instance in the dataset is assigned to the training / validation / test split based on the atomic and compound divergence (similarity) based on weighted distributions using Chernoff coefficient  $C_\alpha(P||Q) = \sum_k p_k^\alpha q_k^{1-\alpha} \in [0, 1]$  (Chung et al., 1989). To make atom distributions similar in train and test, we use  $\alpha = 0.5$  for atom divergence. Here, we set  $\alpha = 0.1$  to reflect that it is more important for a compound to be found in  $P$  (train) rather than the probabilities in  $P$  (train) and  $Q$  (test) match exactly. Following this logic, we define compound divergence, and atom divergence

for a train set  $U$  and test set  $W$  as follows:

$$\begin{aligned} \mathcal{D}_C(U||W) &= 1 - C_{0.1}(\mathcal{F}_C(U) || \mathcal{F}_C(W)) \\ \mathcal{D}_A(U||W) &= 1 - C_{0.5}(\mathcal{F}_A(U) || \mathcal{F}_A(W)) \end{aligned}$$

where  $\mathcal{F}_A(T)$  denotes frequency distribution of atoms, and  $\mathcal{F}_C(T)$  denotes the distribution of compounds for a given set  $T$  and  $D_A$  and  $D_C$  denote atom and compound divergences, respectively. We calculated divergence scores for each instance until the atomic divergence of train and test set  $D_A < 0.02$  and compound divergence of train and test set  $D_C > 0.6$ , which represents a sweet spot in terms of target distributions of atoms and compounds in the train and test sets (see Fig. 4 in the Sec. A.1). Finally, we randomly divide this test set into a validation and test set. The resulting dataset has 8,766 instances, which are split into 4,407 training, 2,184 validation, and 2,175 test instances (see Sec. A for the implementation details and Sec. B for a more detailed analysis of the COMPACT dataset).

## 3 Multimodal Sequential Compositional Generalization

Anticipating what comes next is a fundamental aspect of human cognition (Bar, 2007; Clark, 2015). From a cognitive perspective, it also serves as an engaging training paradigm (Baroni, 2020). In Multimodal Sequential Compositional Generalization, we seek to understand the extent to which multimodal foundation models are capable of understanding what comes next in activity sequences. We propose two tasks to measure multimodal sequential compositional generalization in the COMPACT dataset: (i) next utterance prediction, and (ii) atom classification.

### 3.1 Next Utterance Prediction Task

The next utterance prediction task is a language generation problem, in which models need to predict the text narration that describes the final input in a sequence. Let  $\mathcal{S} = (\mathbf{X}, \mathbf{V}, \mathbf{A})$  denote a triplet representing a short video clip with  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^K$  being a sequence of  $K$  utterances, which describe a household activity and grounded with visual and audio signals, denoted by  $\mathbf{V} = \{\mathbf{v}_i\}_{i=1}^K$  and  $\mathbf{A} = \{\mathbf{a}_i\}_{i=1}^K$ , respectively. This task involves generating the  $(K + 1)^{th}$  utterance,  $\mathbf{y} = \mathbf{x}_{K+1}$ , following the preceding  $K$  utterances and multimodal cues. The training data for this task consists of a set of sequences of microsegments,  $\{(\mathcal{S}, \mathbf{y})\}$ .

### 3.2 Atom Classification Task

The atom classification is a simplified form of the next utterance prediction task. Here, a model only needs to predict the verb and noun in the final input, in isolation from generating grammatically correct sentences. As such, it can be approached as a multi-class classification problem. Diverging from conventional action anticipation tasks (Damen et al., 2022; Gammulle et al., 2019; Ke et al., 2019), our unique setup allows us to approach atom classification through a compositional lens, enabling the prediction of verbs and nouns separately. More formally, let  $S = (\mathbf{X}, \mathbf{V}, \mathbf{A})$  denote a triplet representing a video clip with  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^K$  representing a sequence of  $K$  utterances, which describe a household activity grounded with visual and audio signals, denoted by  $\mathbf{V} = \{\mathbf{v}_i\}_{i=1}^K$  and  $\mathbf{A} = \{\mathbf{a}_i\}_{i=1}^K$ , respectively. More specifically, our atom classification task involves predicting the verb or noun in the  $(K + 1)^{th}$  utterance,  $\mathbf{y} = \mathbf{x}_{K+1}^C$ , following the preceding  $K$  utterances and multimodal cues where  $C$  denotes the verb or noun class.

## 4 Models

In our experiments, we benchmark a variety of neural network models on the proposed next utterance prediction and atom classification tasks, including several text-only (unimodal) and multimodal models for better understanding the importance of different modalities in compositional generalization.

### 4.1 Text-only Unimodal Baseline (L)

Our first baseline is a text-only model to account for unexpected biases in COMPACT (Thomason et al., 2019). This is an encoder-decoder Transformer (Vaswani et al., 2017) with a hidden size of 256 units, where each microsegment is encoded within its context. The model is trained using only the textual utterances  $\mathbf{x}_{1:K}$  from the microsegment as the input, and the next utterance  $\mathbf{x}_{K+1}$  as the target, *i.e.* to predict  $p(\mathbf{x}_{K+1}|\mathbf{x}_{1:K})$ . We use the same backbone in all of our multimodal baselines.

### 4.2 Multimodal Baselines

**Vision and Language (VL):** Our Vision and Language baseline encodes both textual and visual context for the next utterance prediction task. This model encodes the textual utterances  $\mathbf{x}_{1:K}$  of each action from microsegments and the keyframe images  $\mathbf{v}_{1:K}$  to predict the next utterance  $\mathbf{x}_{K+1}$ , *i.e.*  $p(\mathbf{y} = \mathbf{x}_{K+1}|\mathbf{x}_{1:K}, \mathbf{v}_{1:K})$ . This model is adapted

from a model that parses a visual scene and learns cross-modal self-attention (Tsai et al., 2019) over textual inputs and visual data. The visual inputs are encoded using pretrained CNN, and the textual inputs are encoded using a Transformer. More specifically, for the visual modality, we extract two types of features: one type represents global visual features, and the other represents object-level features. For the global features, we use a pretrained ResNet50 model (He et al., 2016) with ImageNet weights (Russakovsky et al., 2015). Object-level features are extracted using a Faster-RCNN object detector (Ren et al., 2017) with a ResNet-101 backbone (He et al., 2016) which is pretrained on MSCOCO (Lin et al., 2014) and finetuned on EK-100. We extract visual features from 5 objects for each keyframe. The resulting representation of a visual keyframe is the concatenation of the global and the object-level features. This concatenated vector is projected into a lower-dimensional space with a linear layer. The textual inputs are encoded using a Transformer with a 256D hidden layer. The visual and textual modalities are then encoded by a cross-modal (CM) self-attention mechanism. In this model, we consider two modalities  $\alpha$  and  $\beta$ , sequences of each modalities are denoted as  $X_\alpha \in \mathbb{R}^{T_\alpha \times d_\alpha}$  and  $X_\beta \in \mathbb{R}^{T_\beta \times d_\beta}$ , respectively and  $T_{(\cdot)}$  denotes sequence length and  $d_{(\cdot)}$  denotes feature dimension. In this model,  $\alpha$  is the language modality, and  $\beta$  is the visual modality. In the cross-modal attention, the textual features are the *keys*, and the visual features are the *queries* and *values*, for aligning visual features to textual features. Let the Query be defined as  $Q_\alpha = X_\alpha W_{Q_\alpha}$ , the Keys as  $K_\beta = X_\beta W_{K_\beta}$ , and the Values as  $V_\beta = X_\beta W_{V_\beta}$ , where  $W_{Q_\alpha} \in \mathbb{R}^{d_\alpha \times d_k}$ ,  $W_{K_\beta} \in \mathbb{R}^{d_\beta \times d_k}$  and  $W_{V_\beta} \in \mathbb{R}^{d_\beta \times d_v}$  are learnable weights. The cross-modal self-attention from  $\beta$  to  $\alpha$  is formulated as a latent adaptation  $Y_\alpha \in \mathbb{R}^{T_\alpha \times d_v}$ :

$$Y_\alpha = \text{CM}_{\beta \rightarrow \alpha}(X_\alpha, X_\beta) = \text{softmax} \left( \frac{Q_\alpha K_\beta^\top}{\sqrt{d_k}} \right) V_\beta \quad (1)$$

The output  $Y_\alpha$  has the same length as  $Q_\alpha$ , but it is represented in the feature space of  $V_\beta$ . This enables the model to fuse different modalities, learning an alignment between the visual and textual features (see Eq.1). There are different strategies proposed in the literature for modeling cross-modal interactions and fusing different modalities (Xu et al., 2023). In our vision and language baseline, we fuse different modalities via a self-attention layer over the aligned vision and language features, which are

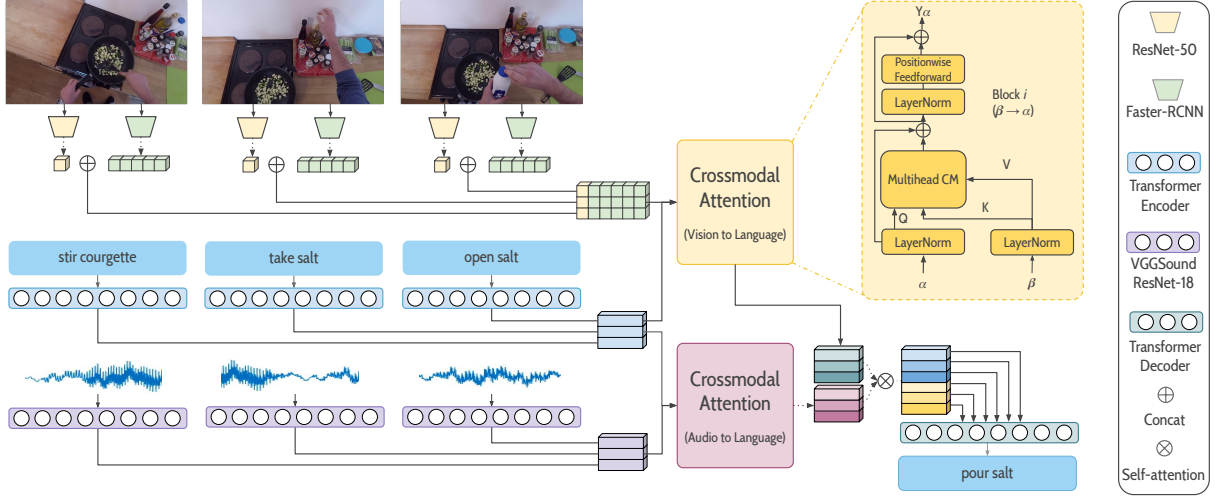


Figure 2: Overview of the AVL baseline which integrates image, object-level, audio, and textual features utilizing two crossmodal attention blocks incorporated within an encoder-decoder Transformer to predict the next utterance.

then fed to a 3-layer Transformer decoder with 4 attention heads that generate the next utterances.

**Audio and Language (AL):** The Audio and Language baseline has the same structure as the VL baseline. The key difference is that we represent the additional context using audio features instead of visual features. The model encodes both the textual utterances  $\mathbf{x}_{1:K}$  and the accompanying audio data  $\mathbf{a}_{1:K}$  to predict the next utterance  $\mathbf{x}_{K+1}$ , *i.e.*  $p(\mathbf{x}_{K+1}|\mathbf{x}_{1:K}, \mathbf{a}_{1:K})$ . The audio features are 512D vectors extracted using VGGSound (Chen et al., 2020), which is pretrained on 200K videos from YouTube, totaling 550 hours of audio data. Here, the proposed AL model learns cross-modal attention over audio and textual features, analogously to the VL model, as inputs to a Transformer decoder.

**Object and Language (OL):** The Object and Language baseline once again uses the same architecture as VL baseline, but the visual context is represented using the labels of detected objects instead of continuous visual features. In this model, we embed object tags as a secondary set of textual features to our model along with the input utterances. Here, the object tags are represented as 292D one-hot encoded vectors (based on the number of unique tags) and projected to 256D with a linear layer. In this case, the cross-modal attention aligns the object tag features with the language features.

**Audio, Vision, and Language (AVL):** In the AVL baseline, we leverage the audio, visual, and textual data using two cross-modal self-attention blocks. We use textual utterances  $\mathbf{x}_{1:K}$  of each

action along with the visual features  $\mathbf{v}_{1:K}$  from the keyframes, and the VGGSound audio features  $\mathbf{a}_{1:K}$  to predict the next utterance  $\mathbf{x}_{K+1}$ , *i.e.*  $p(\mathbf{y} = \mathbf{x}_{K+1}|\mathbf{x}_{1:K}, \mathbf{v}_{1:K}, \mathbf{a}_{1:K})$ . In this model, the input to the Transformer decoder is the concatenation of the audio-aligned textual features from the audio-textual cross-modal block with the visual-aligned textual features from the visual-textual cross-modal block (see Fig. 2 for an overview).

**Object, Audio and Language (OAL):** The OAL baseline model adds an extra modality to the OL baseline model to determine whether audio features affect the performance of a model that uses object tags. Here, we incorporate the extracted audio features from each microsegment into the OL model.

### 4.3 Pretrained Models

To comprehend the significance of large-scale pre-training, we conduct an extensive evaluation involving several publicly available models, namely LLaMA2 (Touvron et al., 2023), IDEFICS (Laurençon et al., 2023), MERLOT Reserve (Zellers et al., 2022, MerlotR), and ImageBind (Girdhar et al., 2023). In our setup, we investigate the performance of encoder-only models across both tasks, whereas auto-regressive models are evaluated exclusively through prompting within the context of the next utterance prediction task. It is worth noting that interpreting the performance of the pretrained models can be complicated as they may violate the distributional consistency between the train and test splits during their pretraining (Kim et al., 2022).

**Unimodal Models:** LLaMA2 is a text-only pre-trained large language model trained on 500B tokens. We evaluate the LLaMA2-Chat 6.7B variant, as this version results in more coherent and relevant predictions due to instruction tuning and RLHF.

**Multimodal Models:** MerlotR learns to extract representations over video frames, text, and audio. The model is composed of an image encoder, an audio encoder, and a joint encoder that fuses textual, visual, and audio representations. This model employs contrastive span training, where an aligned span of audio and text is masked. In its training setup, the objective is to maximize representation similarity to an independent encoding of the masked audio and text spans. We extract multimodal audio and vision features through its pretrained encoder utilizing a similar backbone as in the VL model. ImageBind is a multimodal model that learns joint embeddings for 6 different modalities, including language, vision, and audio. It is trained only on image-paired data to bind the modalities together. We train a decoder using features extracted from the vision, language, and audio modalities. IDEFICS is a large-scale multimodal large language model based on Flamingo (Alayrac et al., 2022) architecture. It is composed of a frozen language model and a frozen vision encoder with learnable cross-attention blocks connecting language and vision modalities. Considering that Flamingo is not publicly available and IDEFICS performs better than other open-source Flamingo implementations such as OpenFlamingo (Awadalla et al., 2023), we experiment with IDEFICS 9B version as the vision LLM. We prompt these models without any finetuning and report 5-shot results for LLaMA2 and IDEFICS (see Sec. A.4 for prompting formats and Sec. E.3 for prompting ablations).

#### 4.4 Task-Specific Changes

For atom classification, we modify the previously described models by altering their architectures. Specifically, we replace the decoder Transformer with two fully connected layers. We train these models with a classification objective conditioned on predicting verbs or nouns in atom classification.

## 5 Experimental Setup

**Evaluation Metrics:** We use unigram BLEU (Papineni et al., 2002), Exact Match (EM), Categorical Accuracy (CA) and BERTScore (Zhang et al., 2019) metrics. The reported values represent the

mean and standard deviation across 3 separate runs. In LLaMA2 and IDEFICS, we use nucleus sampling instead of separate runs. For EM, we calculate the accuracy between the generated text sequence and the ground truth. CA uses the verb and noun categories in EK-100 and calculates the accuracy based on category match between the prediction and ground truth, e.g. the verbs *slice*, *dice*, and *chop* fall into the same verb category *cut*, and the nouns CHEDDAR, PANEER and PARMESAN are grouped into the same noun category CHEESE. Therefore, the *slice* PANEER prediction is deemed accurate if the ground truth is *dice* PARMESAN.

**Training Procedure:** In the next utterance prediction task, models are trained to minimize the negative log-likelihood of generating the next utterance, where the multimodal models are conditioned on additional modalities. Given microsegment  $\mathcal{S}$  and model parameters  $\theta$ , the objective function is to minimize the negative log-likelihood of the  $m$  tokens in the next utterance:  $L(\theta) = -\sum_{i=1}^m \log p(y_i | \mathcal{S}; \theta)$ . In the atom classification task, models are trained by attaching an MLP with a multi-class classification layer to the encoding of a microsegment  $\mathcal{S}$ . The objective function is to minimize the cross-entropy loss of predicting the expected atom (verb or noun):  $L(\theta) = -\log p(\mathbf{x}_{K+1}^C | \mathcal{S}; \theta)$  (see Sec. D for details).

## 6 Results

### 6.1 Next Utterance Prediction

Table 1: Next Utterance Prediction results on the test split. Using audio, visual, or object features consistently improves performance compared to the language-only baseline. The best results are bolded, while the second-best results are underlined. Here, ImageB. and BERTSc. refer to ImageBind and BERTScore, respectively.

Inputs	BLEU	EM	CA	BERTSc.
L	21.75 $\pm$ 1.0	2.89 $\pm$ 0.3	6.43 $\pm$ 0.2	79.06 $\pm$ 0.1
VL	31.25 $\pm$ 0.3	7.27 $\pm$ 0.1	12.95 $\pm$ 0.4	81.27 $\pm$ 0.1
AL	30.82 $\pm$ 0.5	6.81 $\pm$ 0.5	<u>13.22 <math>\pm</math> 0.9</u>	81.20 $\pm$ 0.0
AVL	31.73 $\pm$ 0.4	7.04 $\pm$ 0.4	12.93 $\pm$ 0.8	81.50 $\pm$ 0.1
OL	30.79 $\pm$ 0.6	6.36 $\pm$ 0.2	12.21 $\pm$ 0.1	81.23 $\pm$ 0.1
OAL	<u>32.02 <math>\pm</math> 0.2</u>	<u>7.32 <math>\pm</math> 0.6</u>	13.08 $\pm$ 0.9	<u>81.51 <math>\pm</math> 0.1</u>
MerlotR	31.50 $\pm$ 0.3	6.75 $\pm$ 0.2	12.85 $\pm$ 0.1	81.37 $\pm$ 0.2
ImageB.	<b>33.52 <math>\pm</math> 0.3</b>	<b>9.45 <math>\pm</math> 0.5</b>	<b>15.04 <math>\pm</math> 1.0</b>	<b>82.31 <math>\pm</math> 0.2</b>
IDEFICS	25.64 $\pm$ 0.4	5.76 $\pm$ 0.1	7.89 $\pm$ 0.5	80.92 $\pm$ 0.1
LLaMA2	27.50 $\pm$ 0.6	5.36 $\pm$ 0.6	7.41 $\pm$ 0.7	78.76 $\pm$ 0.2

In Table 1, we present the results of the next utterance prediction experiments. Notably, all mul-










Inputs (utterances and auxiliary modalities)			Prediction (next utterance)	
			GT	: place bowl
clean bowl	open dishwasher	open drawer	L	: close dishwasher
			OL	: place bowl
			VL	: close dishwasher
			AL	: close drawer
			AVL	: dry bowl
			OAL	: place bowl
			MerlotR	: close dishwasher
			ImageBind	: put bowl in dishwasher
			LLaMA2	: get clean
			IDEFICS	: put bowl away.
			GT	: sponge mug
put pan in drainer	pick_up mug	pick_up sponge	L	: put sponge
			OL	: sponge mug
			VL	: sponge mug
			AL	: sponge mug
			AVL	: sponge mug
			OAL	: sponge mug
			MerlotR	: sponge mug
			ImageBind	: sponge mug
			LLaMA2	: put sponge in sink
			IDEFICS	: sponge mug
			GT	: cut pepper
put_down spring_onions	take courgettes	take pepper	L	: cut spring_onions
			OL	: put_down courgette
			VL	: put_down spring_onions
			AL	: put pepper
			AVL	: cut pepper
			OAL	: put_down courgette
			MerlotR	: open pepper
			ImageBind	: cut pepper
			LLaMA2	: take spring onions
			IDEFICS	: put down knife

Figure 3: Next Utterance Prediction qualitative results. Models consider different combinations of input modality, as described in Section 4. In the predictions, **blue** refers to correct, **orange** incorrect and **purple** semantically close.

timodal models surpass the language-only baseline. Our baseline model that incorporates visual features (VL) exhibits consistent increases, showing gains of up to 9 BLEU, 4 EM, 6 CA, and 1 BERTScore points, compared to the language-only variant. Furthermore, harnessing a mix of audio, visual, and language features (AVL) or augmenting audio features with object tags (OAL) leads to additional improvements, emphasizing the contribution of fusing multiple modalities. The most significant boost in performance is observed when visual features are utilized in the ImageBind pre-trained model, resulting in approximate increases of 11, 6, 8, and 2 points, for BLEU, EM, CA, and BERTScore metrics, respectively. The fact that LLaMA2 generates utterances with higher BLEU but lower BERTScore than IDEFICS might suggest that LLaMA2 better imitates the required vocabulary than IDEFICS, even though IDEFICS produces more semantically plausible outputs. Conclusively, ImageBind’s performance shows the advantages of employing pretrained multimodal features over merely merging separate unimodal encodings. In Fig. 3, we present a qualitative comparison of the baseline models via randomly selected examples from the test set. We believe that these illustrative examples effectively showcase the intricate and challenging nature of the proposed COMPACT dataset. During training, the models have never

encountered compounds like *place BOWL*, or *cut PEPPER*. In all of these illustrative examples, the text-only unimodal model fails to generalize to these novel compositions. However, in the first example, the OL and OAL baselines can predict the target composition correctly. ImageBind and IDEFICS, even though not exact matches, generate semantically plausible predictions. In the third example, all multimodal models correctly predict the next utterance by leveraging the auxiliary modalities. Note that, LLaMa2 also fails in this example whereas IDEFICS can generate correct utterances. In the fourth example, the AVL model and ImageBind models correctly predict the *cut PEPPER* utterances. Interestingly, for this example both audio and vision inputs are needed, indicating that for sequential compositional generalization, models might have to leverage the available signal coming from different modalities at the same time.

## 6.2 Atom Classification

In Table 2, we present the outcomes of our atom classification task, which seeks to understand models’ abilities to predict verb and noun atoms in isolation.

For predicting verbs, we observe a similar trend in performance with the next utterance prediction results. However, all models perform poorly in predicting nouns compared to the MRH baseline

Table 2: Quantitative results for Atom Classification. The best and the second-best performing results are highlighted in bold and underlined, respectively.

		EM	CA	BERTScore
Verb Classification	L	13.37 $\pm$ 0.5	28.47 $\pm$ 3.1	75.16 $\pm$ 0.6
	VL	14.02 $\pm$ 0.2	28.68 $\pm$ 2.3	75.29 $\pm$ 0.5
	AL	13.76 $\pm$ 0.3	30.05 $\pm$ 4.6	<u>76.26 <math>\pm</math> 0.5</u>
	AVL	<u>14.06 <math>\pm</math> 0.7</u>	30.98 $\pm$ 2.1	76.12 $\pm$ 0.7
	OL	12.79 $\pm$ 0.1	29.97 $\pm$ 1.5	75.66 $\pm$ 0.2
	OAL	13.91 $\pm$ 0.5	29.90 $\pm$ 1.3	75.97 $\pm$ 0.9
	MerlotR	13.71 $\pm$ 0.1	<b>33.50 <math>\pm</math> 1.8</b>	76.07 $\pm$ 0.2
	ImageBind	<b>15.40 <math>\pm</math> 0.2</b>	<u>31.54 <math>\pm</math> 2.5</u>	<b>76.54 <math>\pm</math> 0.4</b>
MRH	2.39	9.61	73.60	
Noun Classification	L	44.91 $\pm$ 0.3	51.83 $\pm$ 0.3	86.27 $\pm$ 0.2
	VL	42.72 $\pm$ 0.7	49.57 $\pm$ 0.3	85.81 $\pm$ 0.2
	AL	43.95 $\pm$ 0.2	51.11 $\pm$ 0.5	86.08 $\pm$ 0.1
	AVL	43.34 $\pm$ 0.4	50.43 $\pm$ 0.9	85.92 $\pm$ 0.1
	OL	44.35 $\pm$ 0.9	51.24 $\pm$ 0.8	86.00 $\pm$ 0.3
	OAL	43.83 $\pm$ 0.9	51.03 $\pm$ 0.5	85.92 $\pm$ 0.2
	MerlotR	<u>45.42 <math>\pm</math> 0.6</u>	<u>52.24 <math>\pm</math> 0.7</u>	<u>86.42 <math>\pm</math> 0.1</u>
	ImageBind	33.67 $\pm$ 0.3	44.55 $\pm$ 0.1	83.96 $\pm$ 0.1
MRH	<b>57.24</b>	<b>61.15</b>	<b>89.75</b>	

(Most Recent Heuristic). This baseline employs the most recently referenced object in the input microsegment as a prediction for the target noun, and the most recently referenced verb as a prediction for the target verb. While the language-only baseline outperforms the multimodal baselines in noun prediction, we observe an improvement over the language-only model in predicting verbs within the multimodal models. This subtle, yet noteworthy improvement underlines the value of leveraging multiple modalities for verb-related predictions.

### 6.3 Random Split Experiments

The COMPACT dataset, designed to highlight out-of-distribution characteristics, features distinct compositional distributions between training and testing splits, as delineated in Figure 4. This intentional design underscores the out-of-distribution characteristics essential for evaluating the models’ generalization abilities. We expand our analysis to include both in-domain and out-of-domain data, enabling a more comprehensive evaluation of model capabilities and their generalization potential. This analysis aims to compare baseline models on two distinct datasets: our original dataset, which emphasizes compositionality (out-domain), and a new, non-compositional dataset (in-domain) created through random splits of the EK-100 dataset.

Table 3: In-domain Next Utterance Prediction Results

	BLEU	EM	CA	BERTScore
L	23.37	8.95	12.01	80.09
VL	36.49	18.74	22.92	82.83
AL	37.04	19.81	23.86	83.07
AVL	36.08	17.88	21.56	82.84
OL	36.62	19.87	24.22	83.11
OAL	38.59	21.86	25.54	83.64
MerlotR	38.53	21.40	25.82	83.64
ImageBind	40.86	22.62	26.86	84.28
IDEFICS	35.09	17.43	19.11	83.61
LLaMA	30.20	12.10	14.09	79.90

Table 4: In-domain Verb Classification Results

	EM	CA	BERTScore
L	20.26	33.90	77.65
VL	21.90	33.19	77.65
AL	21.44	33.73	77.88
AVL	22.19	34.57	78.17
OL	20.04	33.77	77.51
OAL	21.39	33.21	77.55
MerlotR	21.50	34.39	77.86
ImageBind	21.48	34.00	77.77

Table 5: In-domain Noun Classification Results

	EM	CA	BERTScore
L	48.42	53.45	86.83
VL	49.08	54.51	87.11
AL	47.75	53.67	86.76
AVL	47.47	53.05	86.57
OL	48.55	54.02	86.86
OAL	48.99	54.20	87.02
MerlotR	48.07	53.30	86.72
ImageBind	44.45	52.19	86.26

Our in-depth analysis, reflected in Table 3, indicates a marked improvement in model performance when dealing with in-domain data.

We extend our analysis for the next utterance prediction task in Table 3, to the atom classification task (see Table 4 and Table 5). The results from these analyses consistently show a significant performance improvement for models in the in-domain (non-compositional) setup compared to the compositionally challenging split. This highlights the added complexity and challenge introduced by compositionality, which is not present in

standard atomic classification tasks. The results specifically illustrate that, despite models being trained with examples that include every individual primitive (nouns and verbs), their performance decline when tasked with generalizing to novel compositions. This picture does not change even when most similar examples are used as demonstrations within an in-context learning setup. This decline in performance is not merely due to unfamiliarity with certain primitives but stems from the inherent challenge of understanding and predicting new combinations of these primitives. This divergence in performance between in-domain (random split) and out-of-domain data demonstrates the unique contribution of our work. It underscores the challenge we introduce: pushing the boundaries of current models to not only recognize but also effectively generalize and adapt to novel compositional structures.

## 7 Related Work

**Compositionality.** Baroni (2020) studied the linguistic generalization capabilities of artificial neural networks. Lake et al. (2019) explored compositionality in a human-like few-shot setting, while others have studied compositionality at the representation level such as (Dasgupta et al., 2018; Ettinger et al., 2018). Unimodal compositional generalization datasets such as SCAN (Lake and Baroni, 2017), CFQ (Keysers et al., 2020), and COGS (Kim and Linzen, 2020) have been widely used in the literature to assess generalization abilities of neural networks. In parallel, researchers have been exploring different directions towards compositional generalization *e.g.* meta-learning, (Lake, 2019), altering existing architectures (Akyurek and Andreas, 2021), and data augmentation (Qiu et al., 2022a).

**Grounded Learning.** Zhou et al. (2023) studied grounded learning in a multimodal procedural setup. Wu et al. (2022) benchmarked reasoning and sequencing capabilities of models in a grounded multimodal instructional setting. Johnson et al. (2017) studied systematic generalization in visual reasoning tasks. Bahdanau et al. (2019) investigated systematic generalization in a VQA-like context while Nikolaus et al. (2019) focused on compositionality to construct unseen combinations of concepts while describing images. Seo et al. (2020) transcribed speech to rank correct utterances in instructional videos. Surís et al. (2020) studied compositionality in word acquisition from

narrated videos. Jin et al. (2020) investigated continual learning in unseen compound acquisition from paired image-caption streams.

Other existing studies revolve around crafting conceptual benchmark datasets specifically designed to evaluate compositionality, *e.g.* (de Vries et al., 2019; Vani et al., 2021). Grounded compositional generalization is explored in (Ruis et al., 2020; Wu et al., 2021) within a 2D grid environment. Xu et al. (2021); Yun et al. (2023) investigated grounded compositional generalization for the concept learning problem. Li et al. (2022a) studied compositionality in a grounded setup with audio-language, while Chen et al. (2021) leveraged audio-vision modality pairs.

**Foundation Models.** Recently, researchers have been studying foundation models to explore the possibilities of utilizing different modalities such as audio, vision, and text to solve grounded real-world problems (Guzhov et al., 2022; Girdhar et al., 2023; Driess et al., 2023). More recently, to assess visually-grounded compositional generalization capabilities of models, Bogin et al. (2021) proposed COVR, Zhuo et al. (2023) proposed ViLPAct, Ma et al. (2023) proposed CREPE. Unlike previous studies, in this work, we focus on a real-world audio, vision, and language setting for compositional generalization (see Fig. 1 for an overview). We believe this study contributes toward a better understanding of the open challenges in multimodal sequential compositional generalization for foundation models and spur interest in this direction.

## 8 Conclusion

In this paper, we investigate linguistic compositionality and systematic generalization in a grounded setting for multimodal sequential compositional generalization. We show how a multimodal dataset can be utilized as a challenging test bed for this purpose. We design next utterance prediction and atom classification tasks adopting a methodical approach in generating the training, validation and test sets for our compositional splits. We experiment with several baseline models and investigate models’ ability to generalize to novel compositions and show how multimodal data can contribute towards solving systematic generalization problem and highlight major challenges. We hope our work will stimulate further research in these directions.

## 9 Limitations

Despite the promising results, there are a few limitations of our work. In our work, we introduce a novel dataset called COMPACT carefully curated from the EK-100 dataset (Damen et al., 2022), which involves videos of daily kitchen activities, to dissect the impact of visual and auditory signals on linguistic compositionality. Hence, our conclusions may hinge upon certain domain-specific variables. It could be interesting to conduct future studies in an open-domain setting which might unravel additional insights *e.g.* (Grauman et al., 2022). We investigate several different multimodal models for both the next utterance prediction and atom classification tasks. However, it is important to note that for multimodal learning how to integrate different modalities is considered as an open research problem. In the literature, different strategies for multimodal data fusion have been proposed. Our experimental analysis could be further extended by considering some models that fuse the modalities in a way different than ours. More interestingly, from a systematic generalization point of view, an analysis could be carried out to explore the most effective fusion scheme. Finally, we acknowledge the textual utterances that we use in our work are inherently simplistic and do not capture all of the complexities in natural languages. Consequently, extending this work to a more natural source of language data that mirrors those complexities could be quite interesting direction for future research.

## Ethics Statement

We curated our COMPACT dataset using the video clips from the published EPIC-Kitchens-100 dataset (Damen et al., 2022), which is publicly available under CC BY-NC 4.0 DEED license. The videos that exist in this dataset were recorded voluntarily by the participants who were not financially rewarded. To the best of our knowledge, the proposed systems or models do not pose any risk in terms of fairness, environmental impact and unintended harm or bias.

## Acknowledgment

This work was partly supported by the KUIS AI Center Fellowship to Osman Batur İnce. This work was supported by a research grant (VIL53122) from VILLUM FONDEN.

## References

- Ekin Akyurek and Jacob Andreas. 2021. [Lexicon learning for few shot sequence modeling](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4934–4946, Online. Association for Computational Linguistics.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Yitzhak Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. [Openflamingo: An open-source framework for training large autoregressive vision-language models](#). *ArXiv*, abs/2308.01390.
- Dzmitry Bahdanau, Shikhar Murty, Michael Noukhovitch, Thien Huu Nguyen, Harm de Vries, and Aaron C. Courville. 2019. Systematic generalization: What is required and can it be learned? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Moshe Bar. 2007. The proactive brain: using analogies and associations to generate predictions. *Trends in cognitive sciences*, 11(7):280–289.
- Marco Baroni. 2020. Linguistic generalization and compositionality in modern artificial neural networks. *Phil. Trans. R. Soc. B*, 375(1791):20190307.
- Ben Bogin, Shivanshu Gupta, Matt Gardner, and Jonathan Berant. 2021. [COVR: A test-bed for visually grounded compositional generalization with real images](#). *CoRR*, abs/2109.10613.
- Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 721–725. IEEE.
- Yanbei Chen, Yongqin Xian, A Koepke, Ying Shan, and Zeynep Akata. 2021. Distilling audio-visual knowledge by compositional contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7016–7025.
- JK Chung, PL Kannappan, CT Ng, and PK Sahoo. 1989. Measures of distance between probability distributions. *Journal of mathematical analysis and applications*, 138(1):280–292.
- Andy Clark. 2015. *Surfing uncertainty: Prediction, action, and the embodied mind*. Oxford University Press.

- Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. 2021. [The devil is in the detail: Simple tricks improve systematic generalization of transformers](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 619–634. Association for Computational Linguistics.
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. 2022. [Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100](#). *International Journal of Computer Vision (IJCV)*, 130:33–55.
- Ishita Dasgupta, Demi Guo, Andreas Stuhlmüller, Samuel J. Gershman, and Noah D. Goodman. 2018. [Evaluating compositionality in sentence embeddings](#). *CoRR*, abs/1802.04302.
- Harm de Vries, Dzmitry Bahdanau, Shikhar Murty, Aaron C. Courville, and Philippe Beaudoin. 2019. [CLOSURE: assessing systematic generalization of CLEVR models](#). In *Visually Grounded Interaction and Language (ViGIL), NeurIPS 2019 Workshop, Vancouver, Canada, December 13, 2019*.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. 2023. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*.
- Allyson Ettinger, Ahmed Elgohary, Colin Phillips, and Philip Resnik. 2018. [Assessing composition in sentence vector representations](#). In *Proceedings of the 27th International Conference on Computational Linguistics, COLING 2018, Santa Fe, New Mexico, USA, August 20-26, 2018*, pages 1790–1801. Association for Computational Linguistics.
- Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. 2019. Predicting the future: A jointly learnt model for action anticipation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5562–5571.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190.
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. 2022. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012.
- Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. 2022. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 976–980. IEEE.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on CVPR*, pages 770–778.
- Felix Hill, Andrew K. Lampinen, Rosalia Schneider, Stephen Clark, Matthew Botvinick, James L. McClelland, and Adam Santoro. 2019. [Emergent systematic generalization in a situated agent](#). *CoRR*, abs/1910.00571.
- Xisen Jin, Junyi Du, and Xiang Ren. 2020. [Visually grounded continual learning of compositional semantics](#). *CoRR*, abs/2005.00785.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. 2017. [CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1988–1997. IEEE Computer Society.
- Qiuhong Ke, Mario Fritz, and Bernt Schiele. 2019. Time-conditioned action anticipation in one shot. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9925–9934.
- Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashubin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov, Xiao Wang, Marc van Zee, and Olivier Bousquet. 2020. [Measuring compositional generalization: A comprehensive method on realistic data](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Najoung Kim and Tal Linzen. 2020. [COGS: A compositional generalization challenge based on semantic interpretation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9087–9105, Online. Association for Computational Linguistics.
- Najoung Kim, Tal Linzen, and Paul Smolensky. 2022. Uncontrolled lexical exposure leads to overestimation of compositional generalization in pretrained models. *arXiv preprint arXiv:2212.10769*.
- Brenden M. Lake. 2019. [Compositional generalization through meta sequence-to-sequence learning](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 9788–9798.
- Brenden M. Lake and Marco Baroni. 2017. [Still not systematic after all these years: On the compositional skills of sequence-to-sequence recurrent networks](#). *CoRR*, abs/1711.00350.

- Brenden M. Lake and Marco Baroni. 2018. [Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks](#). In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 2879–2888. PMLR.
- Brenden M. Lake, Tal Linzen, and Marco Baroni. 2019. [Human few-shot learning of compositional instructions](#). In *Proceedings of the 41th Annual Meeting of the Cognitive Science Society, CogSci 2019: Creativity + Cognition + Computation, Montreal, Canada, July 24-27, 2019*, pages 611–617. cognitivescience-society.org.
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M Rush, Douwe Kiela, et al. 2023. Obelisc: An open web-scale filtered dataset of interleaved image-text documents. *arXiv preprint arXiv:2306.16527*.
- Angeliki Lazaridou, Adhi Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d’Autume, Tomas Kocisky, Sebastian Ruder, et al. 2021. Mind the gap: Assessing temporal generalization in neural language models. *Advances in Neural Information Processing Systems*, 34:29348–29363.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*.
- Judith Yue Li, Aren Jansen, Qingqing Huang, Ravi Ganti, Joonseok Lee, and Dima Kuzmin. 2022a. Maqa: A multimodal qa benchmark for negation. In *NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research*.
- Junlong Li, Guangyi Chen, Yansong Tang, Jinan Bao, Kun Zhang, Jie Zhou, and Jiwen Lu. 2022b. Gain: On the generalization of instructional action understanding. In *The Eleventh International Conference on Learning Representations*.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. [Microsoft COCO: common objects in context](#). In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pages 740–755. Springer.
- Adam Liska, Tomas Kocisky, Elena Gribovskaya, Tayfun Terzi, Eren Sezener, Devang Agrawal, D’Autume Cyprien De Masson, Tim Scholtes, Manzil Zaheer, Susannah Young, et al. 2022. Streamingqa: A benchmark for adaptation to new knowledge over time in question answering models. In *International Conference on Machine Learning*, pages 13604–13622. PMLR.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. [What makes good in-context examples for GPT-3?](#) In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. 2023. Crepe: Can vision-language foundation models reason compositionally? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10910–10921.
- Mitja Nikolaus, Mostafa Abdou, Matthew Lamm, Rahul Aralikkatte, and Desmond Elliott. 2019. [Compositional generalization in image captioning](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning, CoNLL 2019, Hong Kong, China, November 3-4, 2019*, pages 87–98. Association for Computational Linguistics.
- Made Nindyatama Nityasya, Haryo Akbarianto Wibowo, Alham Fikri Aji, Genta Indra Winata, Radityo Eko Prasajo, Phil Blunsom, and Adhiguna Kuncoro. 2023. On "scientific debt" in nlp: A case for more rigour in language model pre-training research. *arXiv preprint arXiv:2306.02870*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Linlu Qiu, Peter Shaw, Panupong Pasupat, Pawel Nowak, Tal Linzen, Fei Sha, and Kristina Toutanova. 2022a. [Improving compositional generalization with latent structure and data augmentation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4341–4362, Seattle, United States. Association for Computational Linguistics.
- Linlu Qiu, Peter Shaw, Panupong Pasupat, Tianze Shi, Jonathan Herzig, Emily Pitler, Fei Sha, and Kristina Toutanova. 2022b. Evaluating the impact of model scale for compositional generalization in semantic parsing. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9157–9179.
- Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2017. [Faster R-CNN: Towards real-time object detection with region proposal networks](#). *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6):1137–1149.

- Laura Ruis, Jacob Andreas, Marco Baroni, Diane Bouchacourt, and Brenden M. Lake. 2020. [A benchmark for systematic generalization in grounded language understanding](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252.
- Paul Hongsuck Seo, Arsha Nagrani, and Cordelia Schmid. 2020. [Look before you speak: Visually contextualized utterances](#). *CoRR*, abs/2012.05710.
- Dídac Surís, Dave Epstein, Heng Ji, Shih-Fu Chang, and Carl Vondrick. 2020. [Learning to learn words from visual scenes](#). In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXIX*, volume 12374 of *Lecture Notes in Computer Science*, pages 434–452. Springer.
- Jesse Thomason, Daniel Gordon, and Yonatan Bisk. 2019. [Shifting the baseline: Single modality performance on visual navigation & QA](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1977–1983, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248.
- Hugo Touvron, Louis Martin, and Kevin Stone et al. 2023. [Llama 2: Open foundation and fine-tuned chat models](#).
- Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2019. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Florence, Italy. ACL.
- Ankit Vani, Max Schwarzer, Yuchen Lu, Eeshan Dhekane, and Aaron Courville. 2021. [Iterated learning for emergent systematicity in VQA](#). In *International Conference on Learning Representations*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Te-Lin Wu, Alex Spangher, Pegah Alipoormolabashi, Marjorie Freedman, Ralph Weischedel, and Nanyun Peng. 2022. Understanding multimodal procedural knowledge by sequencing multimodal instructional manuals. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4525–4542.
- Zhengxuan Wu, Elisa Kreiss, Desmond C Ong, and Christopher Potts. 2021. Reascan: Compositional reasoning in language grounding. *arXiv preprint arXiv:2109.08994*.
- Guangyue Xu, Parisa Kordjamshidi, and Joyce Chai. 2021. [Zero-shot compositional concept learning](#). In *Proceedings of the 1st Workshop on Meta Learning and Its Applications to Natural Language Processing, pages 19–27*, Online. Association for Computational Linguistics.
- Peng Xu, Xiatian Zhu, and David A Clifton. 2023. Multimodal learning with transformers: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Tian Yun, Usha Bhalla, Ellie Pavlick, and Chen Sun. 2023. [Do vision-language pretrained models learn composable primitive concepts?](#) *Trans. Mach. Learn. Res.*, 2023.
- Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. 2022. Merlot reserve: Neural script knowledge through vision and language and sound. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16375–16387.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Yu Zhou, Sha Li, Manling Li, Xudong Lin, Shih-Fu Chang, Mohit Bansal, and Heng Ji. 2023. Non-sequential graph script induction via multimedia grounding. *arXiv preprint arXiv:2305.17542*.
- Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. 2019. Cross-task weakly supervised learning from instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3537–3545.
- Terry Yue Zhuo, Yaqing Liao, Yuecheng Lei, Lizhen Qu, Gerard de Melo, Xiaojun Chang, Yazhou Ren, and Zenglin Xu. 2023. [ViLPAct: A benchmark for compositional generalization on multimodal human activities](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2192–2207, Dubrovnik, Croatia. Association for Computational Linguistics.

## Appendix

In the following, we provide a comprehensive set of supplementary notes that delve deeper into various aspects of our research:

- **Data Curation, Algorithms, and Preprocessing (Section A):** This section outlines the steps taken in data curation, algorithmic processes, and preprocessing techniques applied.
- **Exploratory Analysis of COMPACT (Section B):** Here, we present a detailed analysis of the COMPACT dataset, highlighting its unique characteristics.
- **Implementation Details and Reproducibility (Section D):** This section offers a detailed account of our implementation methodology, providing valuable information for those interested in replicating or extending our work.
- **Further Analysis (Section E):** We conduct additional analyses, expanding on key findings and offering deeper insights into the compositional generalization phenomenon.
- **Ethics Statement (Section 9):** In this section, we present a comprehensive ethics statement detailing our commitment to ethical research practices throughout the study.

### A Data Curation, Algorithms and Preprocessing

#### A.1 Curating COMPACT: an Overview

In our data curation and preprocessing for COMPACT, we leverage the EPIC-Kitchens-100 (EK-100) dataset, a collection of egocentric kitchen activity videos, which are split into shorter clips with accompanying narrations and audio tracks—referred to as “microsegments”.

To curate our sequences of microsegments, we employ a window of 4 clips, with the initial 3 clips serving as context and the last one designated for prediction, yielding a total of 22,136 instances. We filter out repeated utterances that represent a continuation of the same action, treating them as duplicates.

Additionally, we exclusively consider text descriptions that share common nouns, ensuring that the noun mentioned in the target description also appears in the source text. This heuristic guarantees the presence of the target noun in both input

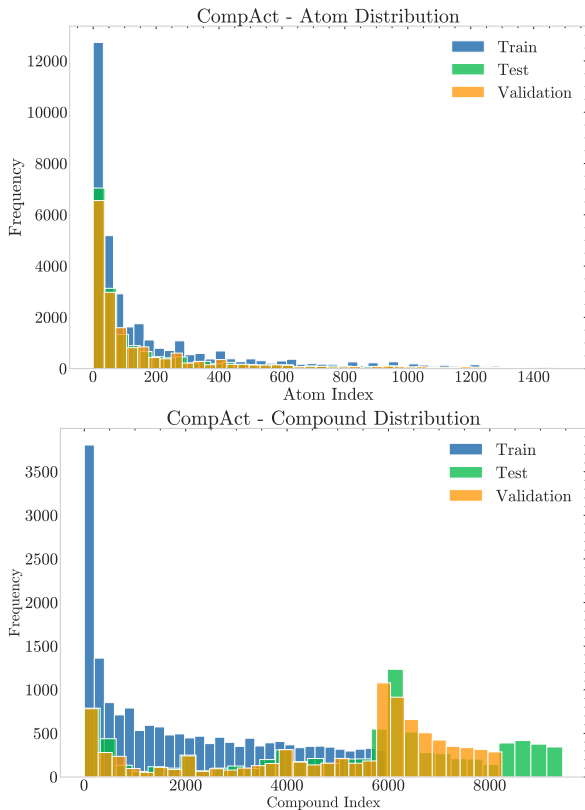


Figure 4: Plot on the top demonstrates the distribution of atoms while the plot on the bottom shows the distribution of compounds for the train/validation/test splits in compositional split setup.

sequences during both inference and training, allowing our setup to solely evaluate compositionality and systematic generalization.

In our experimental setup, we introduce a scenario where a model must have prior exposure to all constituent atoms within a test instance, such as GRAB THE PLATE. WASH CUCUMBER. TAKE KNIFE., and is then tasked with predicting the subsequent utterance, such as SLICE CUCUMBER, during inference. It is important to emphasize that this target composition has never been encountered during the model’s training phase (refer to Fig. 5). This setup allows testing models’ ability to generalize to entirely unobserved compositions, even those with zero probability of occurrence in the training data. To create such dataset splits, we employ the Maximum Compound Divergence (MCD) heuristic, crafting distributions that maintain similarity in the distribution of individual concepts (atoms), while deliberately introducing disparities in the distributions of concept combinations. In our case, we utilize 97 verb classes and 300 noun classes from the EK-100 dataset as the atoms. In particular,







	Inputs (utterances and auxiliary modalities)			Target (next utterance)
Image				GT : put_down knife
Text	take_off gloves	pick_up knife	check butter in bowl	
Image				GT : rinse pan
Text	put_down water_filter	put_down glass	pick_up pan from sink	

Figure 5: Curating dataset instances for compositional generalization. Targets such as *put\_down* KNIFE and *rinse* PAN have never been observed by the learner during the training phase.

each sample is assigned to a specific split based on the atomic and compound divergence (similarity) based on weighted distributions using Chernoff coefficient (Chung et al., 1989). This process yields 8,766 instances, which are further partitioned into 4,407 for training, 2,184 for validation, and 2,175 for testing.

In Fig. 4, we visualize the atomic and compound distributions over the constructed training, validation, and test splits of our proposed compositional setup. Notably, these splits exhibit similar distributions concerning atoms while training and val/test splits do differ in terms of compounds.

## A.2 Atom and Compound Selection

In Algorithm 1, we describe the heuristic we use to create the compositional splits in COMPACT following the Maximum Compound Divergence (Keysers et al., 2020)

## A.3 Preprocessing

### A.3.1 Choosing Keyframes from Videos

We adopt a straightforward yet effective approach to select representative images from each microsegment. We employ a simple heuristic to identify which keyframes to be selected for the span of the video clip. In particular, we run an object detector on the video frames and select the frames containing the highest count of object proposals detected by the object detector. This selection ensures that we capture the most visually informative frame from among the available candidates. In the case of ImageBind, we opt for the middle frame from each narration video.

---

### Algorithm 1: Split Generation Algorithm

---

**Data:** Dataset  $M$

**Result:** Train split  $U$ , Test split  $W$

```

1 Init  $U, W$ ;
2 Init Atom Divergence  $D_A$ , Compound
  Divergence  $D_C$ ;
3 Init  $M_T$  with items in  $M$ ;
4 Init  $i$  to 0;
5 while  $M_T$  is not empty do
6   Randomly choose  $T \in \{U, W\}$  to add
     an item;
7   if  $i = 0$  then
8     Randomly select and remove an
       item  $m$  from  $M_T$ ;
9     Add  $m$  to split  $T$ ;
10  else
11    Calculate  $D_A$  for remaining items if
      added to  $T$ ;
12    Filter items with  $D_A$  below a
      threshold;
13    if no items meet the criteria then
14      Select item with highest  $D_C$  as
        the best candidate;
15    else
16      Calculate  $D_C$  for items if added
        to  $T$ ;
17      Select the item with highest  $D_C$ 
        as the best candidate;
18    Add the best candidate item to split
       $T$ ;
19  Increment  $i$  by 1;

```

---

```

Predict the next narration given 3 sequential previous narrations from a cooking video
put down bowl . move frying pan . pick up spatula => put down spatula
put down bowl . move jar . pick up egg => crack egg
move yoghurt . put down bowl . pick up yogurt => put yoghurt
put down bowl . grab wok . move tap => lather wok
put down bowl . pick up spatula . stir meat pieces with spatula => put down spatula
pick up tins . put down tins . move bowl =>

```

Figure 6: Prompt template utilized for LLaMA2 evaluation.

```

Predict the next action narration given 3 sequential previous actions (image-narration pairs) in a
cooking video.
put down bowl <Image 1> . move frying pan <Image 2> . pick up spatula <Image 3> => put down
spatula
pick up tins <Image 1> . put down tins <Image 2> . move bowl <Image 3> =>

```

Figure 7: Prompt template utilized for IDEFICS evaluation.

### A.3.2 Tokenization

As a preprocessing step, we replace multiword tokens with a single word. For instance, each occurrence of *put-down* is replaced with *put\_down*, and each occurrence of OLIVE OIL is replaced with OLIVE\_OIL. This preprocessing step is not applied to LLaMA2 and IDEFICS, since these models have their vocabulary. Similarly, LLaMA2 and IDEFICS use their tokenizers while other models simply use a whitespace tokenizer.

### A.4 Prompt Format

In this section, we describe the heuristic we employ to formulate the inputs for our evaluation prompts targeting generative models. It is worth noting that the prompting templates for IDEFICS and LLaMA2, though similar, are not interchangeable as IDEFICS can harness both visual and language data.

First, for both LLMs, we include an instruction at the start of the prompt as our language models are instruction-tuned. Then, we enumerate a set of few-shot examples. Finally, we provide the source section at the end of the prompt, leaving the target to be predicted.

#### A.4.1 LLaMA2 Prompt Example

An example LLaMA2 5-shot prompt can be seen in Fig. 6.

#### A.4.2 IDEFICS Prompt Example

An example IDEFICS 1-shot prompt can be seen in the Fig. 7. <Image  $n$ > denotes the image for the  $n^{th}$  narration scene.

## B Exploratory Analysis of COMPACT

In this section, we share an exploratory analysis of the COMPACT. Fig. 8 illustrates the verb and noun distributions in the COMPACT dataset where the validation and test splits are jointly stacked on top of the train split occurrences and displayed in a lighter color.

## C Choice of EK-100 for COMPACT

The EPIC-Kitchens-100 (EK-100) dataset was chosen due to its established reputation in the research community and its densely annotated instructions, offering a rich and diverse dataset. It also has a clear segmentation of instructions, including verb and noun annotations, making it an ideal candidate for curating the COMPACT dataset, allowing us to leverage audio, vision, and text modalities effectively. Previously proposed datasets in the literature such as CrossTask (Zhukov et al., 2019) and GAIN (Li et al., 2022b) also consists of text instructions and multimodal components. Unlike CrossTask, which focuses on cross-task generalization, our study centers on compositional generalization. Similarly, while sharing similarities with GAIN in dataset formulation and the use of instructional videos, COMPACT differs in the description of atomic concepts and the mathematical definition of out-of-distribution (OOD) scenarios. We also conduct further analysis to evaluate whether the proposed benchmarks such as CrossTask or GAIN could be considered for a compositional generalization benchmark. Nevertheless, the lack of proper annotations for atoms and compounds and the num-

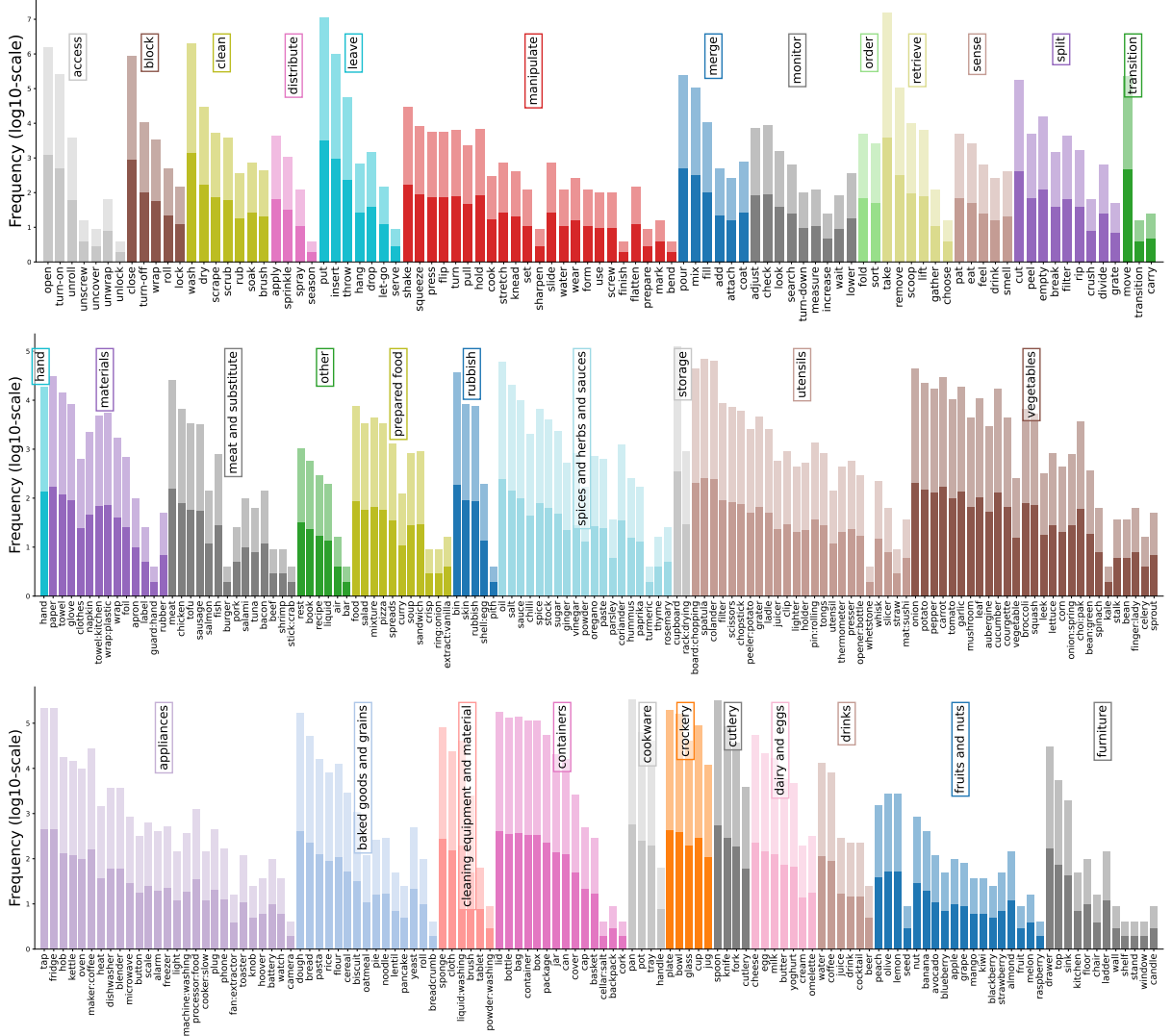


Figure 8: Distribution of verbs (top) and nouns (middle and bottom) from COMPACT

ber of instances seem to be a challenge to generate compositional splits for these benchmarks.

## D Implementation Details and Reproducibility

For the reproducibility of our results, we plan to make the code, models, COMPACT splits, and extracted features publicly available under CC BY-NC 4.0 DEED license. All models are implemented with PyTorch. We use torchtext library for the BLEU metric and evaluate library for the BERTScore metric.

### D.1 Training Regime and Hyperparameters

We use the AdamW optimizer (Loshchilov and Hutter, 2017) with ReduceLROnPlateau learning rate scheduler to reduce the learning rate during training when validation BLEU plateaus. To train the mod-

els for the next utterance prediction, we employ cross-entropy loss, and initialize network weights via uniform distribution for both the encoder and the decoder. We use an early stopping strategy and stop the training if validation BLEU does not improve after a certain threshold. We clip gradients set the gradient threshold to 0.1, and use a 3-layer multihead attention with 4-heads in the crossmodal self-attention block in all our multimodal models. We use the same strategy for atom classification, with one distinction where we use accuracy for early stopping and learning rate scheduler.

For both tasks, we use 50 epochs as the early stopping threshold. We use a dropout rate of 0.3 and the AdamW optimizer with a  $3e-4$  learning rate and  $5e-5$  weight decay. We use the ReduceLROnPlateau learning rate scheduler with patience of 40. Following the insights from Csordás et al. (2021),

Table 6: Model sizes and their training times for our experiments. Training times are averaged over 3 runs.

Model	Next Utterance Prediction		Atom Classification		
	#params	Train Time	#params	Noun Train Time	Verb Train Time
L	4.8M	20:45	2.1M	2:45	2:00
OL	12.0M	38:15	9.3M	12:15	8:30
VL	12.5M	49:15	9.7M	26:45	22:00
AL	12.0M	40:00	9.3M	10:45	8:15
AVL	12.6M	52:00	9.9M	15:30	13:00
OAL	12.1M	39:00	9.4M	14:15	10:30
MerlotR	12.1M	18:30	9.4M	6:30	4:45
ImageBind	8.4M	28:15	5.7M	9:45	8:30

we use the performance score as a monitoring metric for the scheduler (also early stopping) rather than using loss. For the next utterance prediction and action classification tasks, we use BLEU and accuracy scores, respectively.

## D.2 Model Sizes and Training Time

In Table 6, we present the number of trainable parameters and training time (MM:SS) for all of our trainable baseline models for both the next utterance prediction and atom classification tasks.

LLaMA2 and IDEFICS experiments are run on NVIDIA Tesla T4 and NVIDIA Tesla V100 GPUs respectively. Other experiments are run on NVIDIA 1080Ti GPUs.

## E Further Analysis

### E.1 Generalization on Validation Split

Table 7: Next utterance prediction results on validation split. Using audio, visual, or object features always improves performance compared to the language-only unimodal baseline. We report the mean and the standard deviation across three runs.

	BLEU	EM	CA	BERTScore
L	21.43 $\pm$ 0.5	2.88 $\pm$ 0.1	6.22 $\pm$ 0.2	79.20 $\pm$ 0.1
VL	30.59 $\pm$ 0.4	7.35 $\pm$ 0.6	12.39 $\pm$ 1	81.24 $\pm$ 0.4
AL	30.47 $\pm$ 0.1	7.06 $\pm$ 0.3	12.16 $\pm$ 0.2	81.19 $\pm$ 0.1
AVL	31.22 $\pm$ 0.1	7.44 $\pm$ 0.4	12.54 $\pm$ 0.3	81.44 $\pm$ 0.1
OL	30.50 $\pm$ 0.4	7.03 $\pm$ 0.2	12.42 $\pm$ 0.8	81.10 $\pm$ 0.1
OAL	<u>31.42 <math>\pm</math> 0.1</u>	<u>7.99 <math>\pm</math> 0.5</u>	<u>13.36 <math>\pm</math> 0.2</u>	<u>81.50 <math>\pm</math> 0.1</u>
MerlotR	31.36 $\pm$ 0.4	7.17 $\pm$ 0.6	12.68 $\pm$ 0.5	81.34 $\pm$ 0.1
ImageBind	<b>34.13 <math>\pm</math> 0.5</b>	<b>10.45 <math>\pm</math> 0.8</b>	<b>16.08 <math>\pm</math> 0.8</b>	<b>82.45 <math>\pm</math> 0.2</b>
IDEFICS	25.15 $\pm$ 0.8	5.66 $\pm$ 0.5	7.17 $\pm$ 0.5	80.75 $\pm$ 0.2
LLaMA2	26.52 $\pm$ 0.5	5.37 $\pm$ 0.3	6.99 $\pm$ 0.4	78.59 $\pm$ 0.1

In Table 7 we present generalization performance on the validation split for the next utterance prediction task and in Table 8 we demonstrate the generalization performance on validation split for the atom classification task.

Table 8: Atom classification results on validation split. We report the mean across three runs. The best and second best-performing results are highlighted in bold and underlined, respectively.

		EM	CA	BERTScore
Verb Classification	L	12.92 $\pm$ 0.8	28.96 $\pm$ 3.3	74.96 $\pm$ 0.4
	VL	14.02 $\pm$ 0.2	30.23 $\pm$ 1.7	75.19 $\pm$ 0.5
	AL	<u>14.48 <math>\pm</math> 0.7</u>	31.07 $\pm$ 3.7	<u>76.21 <math>\pm</math> 0.2</u>
	AVL	14.01 $\pm$ 0.3	31.15 $\pm$ 2.5	75.84 $\pm$ 0.9
	OL	12.77 $\pm$ 0.3	30.87 $\pm$ 1.1	75.53 $\pm$ 0.3
	OAL	14.30 $\pm$ 0.2	30.79 $\pm$ 1.5	76.02 $\pm$ 0.6
	MerlotR	13.15 $\pm$ 0.5	<b>32.53 <math>\pm</math> 0.6</b>	75.74 $\pm$ 0.2
	ImageBind	<b>14.91 <math>\pm</math> 0.2</b>	<u>31.31 <math>\pm</math> 3.1</u>	<b>76.28 <math>\pm</math> 0.5</b>
Noun Classification	MRH	—	—	—
	L	<u>44.78 <math>\pm</math> 0.8</u>	<u>52.28 <math>\pm</math> 0.8</u>	<u>86.35 <math>\pm</math> 0.1</u>
	VL	42.35 $\pm$ 0.3	49.38 $\pm$ 0.4	85.88 $\pm$ 0.1
	AL	43.71 $\pm$ 0.2	50.79 $\pm$ 0.2	86.05 $\pm$ 0.1
	AVL	43.48 $\pm$ 0.5	50.59 $\pm$ 0.8	86.10 $\pm$ 0.2
	OL	44.03 $\pm$ 0.3	51.40 $\pm$ 0.1	86.03 $\pm$ 0.1
	OAL	44.13 $\pm$ 0.5	50.86 $\pm$ 0.8	86.09 $\pm$ 0.2
	MerlotR	44.52 $\pm$ 0.8	51.31 $\pm$ 0.6	86.22 $\pm$ 0.2
	ImageBind	34.15 $\pm$ 0.5	44.59 $\pm$ 0.2	84.11 $\pm$ 0.1
	MRH	<b>57.51</b>	<b>60.90</b>	<b>89.89</b>

### E.2 Generalization Performance over Epochs

In Fig. 9, we report the BLEU scores of the models over the training, validation, and test splits at different epochs. These plots demonstrate that in a compositional setup, models can perform well in the training set but this does not mean they can generalize to unseen distributions.

### E.3 Prompting Ablations

#### E.3.1 Additional Few-Shot Results

Table 9 and 10 offer interesting insights regarding few-shot compositional capabilities of IDEFICS and LLaMA2 models. First, we see a significant performance discrepancy between IDEFICS and LLaMA2 on zero-shot prediction results. As

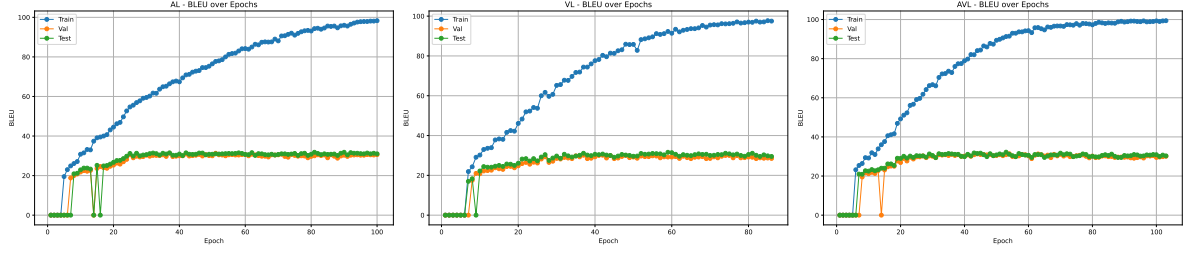


Figure 9: Generalization performance of the AL, VL, and AVL models over the epochs. Even though the training performance of a model improves on COMPACT, this does not necessarily mean that its validation and test performance will also become better due to the compositional nature of the COMPACT dataset.

IDEFICS additionally utilizes visual information over LLaMA2, it displays better zero-shot generalization capabilities. While LLaMA2 outperforms IDEFICS in one-shot and few-shot BLEU scores, contrastingly, IDEFICS outperforms LLaMA2 on BERTScores. As LLaMA2 outperforms the LLM of IDEFICS (instruct-tuned LLaMA1) on many benchmarks (Touvron et al., 2023), we infer that LLaMA2 can imitate the vocabulary of few-shot examples better than IDEFICS, resulting in higher BLEU scores. However, higher BERTScores imply that IDEFICS can reflect the semantics of the ground truth prediction better.

Table 9: Next utterance prediction results on test split for IDEFICS. As the few-shot example count increases, performance improves on every metric consistently.

$k$ -shot	BLEU	EM	CA	BERTScore
0-shot	$8.98 \pm 0.2$	$0.06 \pm 0.1$	$0.12 \pm 0.0$	$75.58 \pm 0.1$
1-shot	$20.25 \pm 0.2$	$4.12 \pm 0.1$	$5.30 \pm 0.1$	$79.75 \pm 0.0$
3-shot	$24.85 \pm 0.7$	$5.37 \pm 0.4$	$7.20 \pm 0.3$	$80.78 \pm 0.1$
5-shot	$25.64 \pm 0.4$	$5.76 \pm 0.1$	$7.89 \pm 0.5$	$80.92 \pm 0.1$
8-shot	$26.18 \pm 0.3$	$6.06 \pm 0.2$	$7.92 \pm 0.3$	$81.19 \pm 0.1$

Table 10: Next utterance prediction results on test split for LLaMA2. As the few-shot example count increases, performance improves on every metric consistently.

$k$ -shot	BLEU	EM	CA	BERTScore
0-shot	$2.02 \pm 3.5$	$0.13 \pm 0.1$	$0.15 \pm 0.1$	$71.68 \pm 0.1$
1-shot	$23.89 \pm 0.7$	$3.98 \pm 0.2$	$5.77 \pm 0.1$	$77.90 \pm 0.2$
3-shot	$26.17 \pm 0.6$	$5.07 \pm 0.1$	$7.00 \pm 0.1$	$78.35 \pm 0.1$
5-shot	$27.50 \pm 0.6$	$5.36 \pm 0.6$	$7.41 \pm 0.7$	$78.76 \pm 0.2$
8-shot	$27.58 \pm 0.3$	$5.60 \pm 0.2$	$8.01 \pm 0.4$	$78.95 \pm 0.1$

### E.3.2 Few-Shot Example Selection

For few-shot example selection, rather than randomly picking  $k$ -shot examples, we employ a simple heuristic. As Liu et al. (2022) highlights that se-

lecting similar examples improves in-context learning performance, we select the most similar  $k$  examples as few-shot examples. The similarity measure between the two examples is based on the noun and verb overlap. First, the intersection between the set of nouns and the set of verbs between the main example and all training examples is computed. If the sum of the cardinality of these sets is the largest between the main example and a few-shot example, the few-shot example is the most similar example of the main example. We provide a validation comparison between the random example selection and our heuristic in Table 11.

Table 11: Next utterance prediction BLEU scores on validation split for IDEFICS for a single run. Greedy decoding is used and the best score is bolded.

Strategy	0-shot	1-shot	3-shot	5-shot
Random selection	<b>28.5</b>	31.0	18.7	19.3
Our heuristic	<b>28.5</b>	<b>34.8</b>	<b>23.4</b>	<b>23.5</b>

### E.3.3 Prompt Template Selection

For IDEFICS, as images should be included in the prompt, the selection of a prompt template is important. We compared two prompt templates (see Fig. 7 and Fig. 10) and after preliminary analysis, used the best-performing template in our paper (see Table 12).

Table 12: Next utterance prediction results on validation split for IDEFICS. Overall, the used template outperforms unused template.

Template	BLEU	EM	CA	BERTScore
Used	<b>25.64 <math>\pm</math> 0.4</b>	$5.76 \pm 0.1$	<b>7.89 <math>\pm</math> 0.5</b>	<b>80.92 <math>\pm</math> 0.1</b>
Unused	$22.19 \pm 0.3$	<b>5.91 <math>\pm</math> 0.4</b>	$7.14 \pm 0.5$	$80.53 \pm 0.1$

Predict the next action narration given 3 sequential previous actions (image-narration pairs) in a cooking video.
Narration 1: put down bowl Image 1: <Image 1>
Narration 2: move frying pan Image 2: <Image 2>
Narration 3: pick up spatula Image 3: <Image 3>
Narration 4: put down spatula
Narration 1: pick up tins Image 1: <Image 1>
Narration 2: put down tins Image 2: <Image 2>
Narration 3: move bowl Image 3: <Image 3>
Narration 4:

Figure 10: The unused alternative prompt template for IDEFICS evaluation.

#### E.4 Unseen Compositions for Pre-trained LLMs

We recognize the challenge of ensuring that pre-trained LLMs have not been exposed to certain compositions during training. Our motivation to explore pre-trained LLMs is inspired by these models’ recent successes in various tasks, and to understand how these models succeed in compositional generalization, we conduct experiments with in-distribution data (randomly generated training/validation/test splits), to implicitly determine the extent of prior exposure to unseen compositions in these models.

Our comprehensive analysis, depicted in Table 3, showcases significant performance improvements in in-domain setting for various baseline models and multimodal LLMs. This contrasts with their performance in out-of-domain setting, highlighting the rigorous nature of our compositional task. The difficulty these models face in generalizing to novel compositions, despite the possibility of exposure to similar examples, indicates a crucial challenge in current multimodal learning.

#### E.5 Video-Based Baseline Experiments

We expand our evaluation to understand the video-based baseline performance through an analysis with the Otter model [Li et al. \(2023\)](#). This instruction-tuned Video Language Model (VLM) processes videos as sequential images and was assessed using few-shot prompting. Our findings, detailed in Table 13, indicate that while Otter performs similarly over OpenFlamingo-9B, the gains with increased context examples are not as substantial as anticipated.

Table 13: Otter Model Results

	BLEU	EM	CA	BERTScore
0-shot	9.31	0.01	0.03	74.54
1-shot	9.26	0.16	0.26	73.30
3-shot	9.95	0.45	0.54	74.21
5-shot	10.5	0.41	0.49	74.83
8-shot	10.85	0.34	0.43	74.38

Table 14: Model size comparison results for OpenFlamingo

	BLEU	EM	CA	BERTScore
OpenFlamingo-3B	8.96	0.09	0.15	73.34
OpenFlamingo-9B	11.15	0.75	0.87	72.38

#### E.6 Comparison of Model Sizes

We conduct additional experiments to compare model performance across different scales using different sizes of the OpenFlamingo model, specifically OpenFlamingo-3B-vitl-mpt1b and OpenFlamingo-9B-vitl-mpt7b. Using a 5-shot prompting approach, we present our evaluation results in Table 14.

Interestingly, the OpenFlamingo-9B’s performance is substantially lower than the IDEFICS results (11 BLEU for OpenFlamingo vs. 24 BLEU for IDEFICS). We attribute this discrepancy primarily to the models’ differing abilities to integrate interleaved images and to the instruction-tuned nature of IDEFICS’s LLM, enhancing its prompt adherence. Although scaling up model size shows some performance improvement, the disparities with other baselines are noteworthy. This observation aligns with ongoing discussions in the field about the impact of model scale on performance, as thoroughly investigated by [Qiu et al. \(2022b\)](#).

#### E.7 Analysis of the Failure Cases

We delve deeper into failure cases, particularly examining instances where a correctly predicted verb is paired with a misclassified noun. Our analysis focus on whether these nouns are more likely to match with training compositions. Additionally, we want to reiterate that our train/val/test split curation in CompAct was meticulously designed to ensure similar primitive coverage across splits while varying compound compositions (see Fig. 4).

In particular, in Table 15, we share the percentage of the misclassified nouns for correctly pre-

Table 15: Percentage of misclassified nouns in train and test split V-N compositions

	Train Ratio	Test Ratio
L	9.97%	7.63%
VL	10.64%	8.04%
AL	10.98%	8.34%
AVL	10.65%	8.21%
OL	10.47%	8.06%
OAL	10.44%	8.12%
MerlotR	10.58%	8.13%
ImageBind	10.39%	8.01%

dicted verbs for the training and test compositions. These ratios are averaged across 3 runs for each model.

### E.8 t-SNE Visualization for Audio and Visual Embeddings

To understand the features we extracted via VG-GSound and ResNet50 backbones and how well they encode the audio and visual spaces, we visualized the raw feature embeddings by projecting them to 2D space via t-SNE. In Fig. 11, we highlight the compounds with the verb *rinse* and observe that audio features can meaningfully encode activities. Similarly, we analyzed the raw global visual embeddings and highlighted the compounds with the noun FRIDGE, the visualization shows the extracted global visual embeddings can effectively encode visual surroundings.

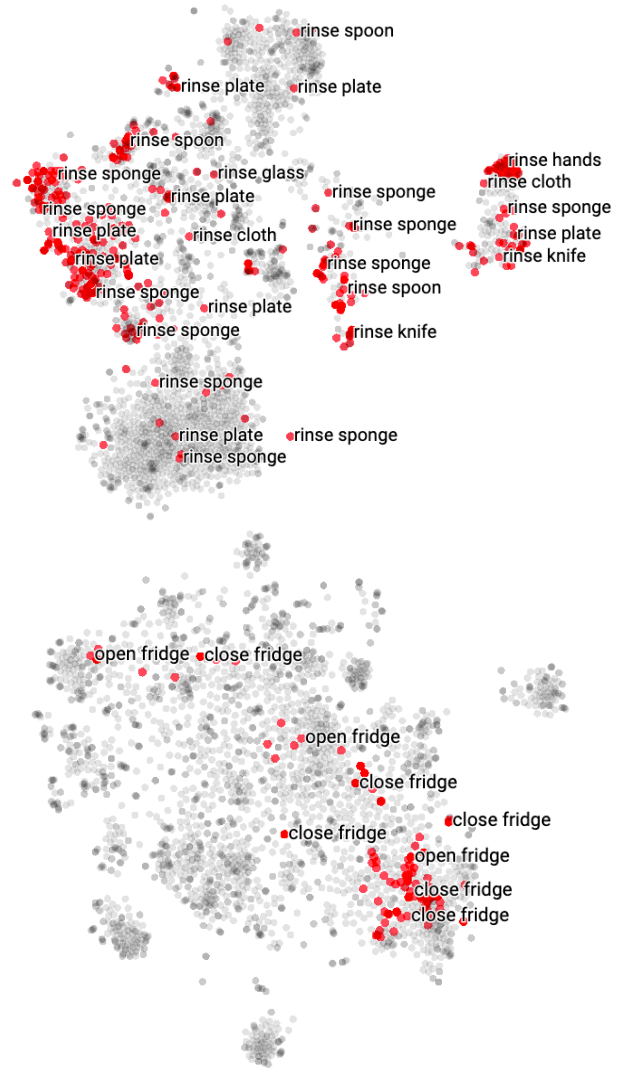


Figure 11: Feature projection to 2D space with t-SNE using raw audio and global visual features. At the top, audio space is shown with the verb *rinse* being specifically highlighted. At the bottom, visual space is given with noun FRIDGE being particularly highlighted. Sampled by most common compounds appearing at least 25 times in COMPACT, equally distributed for each compound ( $N = 25$ ).