# The Curious Case of Representational Alignment: Unravelling Visio-Linguistic Tasks in Emergent Communication

**Tom Kouwenhoven**
LIACS
Leiden University
t.kouwenhoven@liacs.leidenuniv.nl

**Max Peeperkorn**
School of Computing
University of Kent
m.peeperkorn@kent.ac.uk

**Bram van Dijk**
LUMC
Leiden University Medical Center
b.m.a.van_dijk@lumc.nl

**Tessa Verhoef**
LIACS
Leiden University
t.verhoef@liacs.leidenuniv.nl

## Abstract

Natural language has the universal properties of being compositional and grounded in reality. The emergence of linguistic properties is often investigated through simulations of emergent communication in referential games. However, these experiments have yielded mixed results compared to similar experiments addressing linguistic properties of human language. Here we address *representational alignment* as a potential contributing factor to these results. Specifically, we assess the representational alignment between agent image representations and between agent representations and input images. Doing so, we confirm that the emergent language does not appear to encode human-like conceptual visual features, since agent image representations drift away from inputs whilst inter-agent alignment increases. We moreover identify a strong relationship between inter-agent alignment and topographic similarity, a common metric for compositionality, and address its consequences. To address these issues, we introduce an alignment penalty that prevents representational drift but interestingly does not improve performance on a compositional discrimination task. Together, our findings emphasise the key role representational alignment plays in simulations of language emergence.

## 1 Introduction

Human language bears unique properties that make it a powerful tool for communication. A well-known property is compositionality: the ability to combine meaningful words into more complex meanings (Hockett, 1959). The emergence of compositionality is studied extensively in the field of language evolution through human experiments (e.g. Selten and Warglien, 2007; Kirby et al., 2008, 2015; Raviv et al., 2019a). An important finding from this field is that the unique nature of human

language can be explained as a consequence of biases for simplicity and expressivity imposed during continuous language learning and use (Smith, 2022). Computational simulations of language emergence have also been used to study the emergence of linguistic properties (e.g. de Boer, 2006; Steels and Loetzsch, 2012), and have seen a rising interest in the field of computational linguistics (Lazaridou and Baroni, 2020). Here, compositionality in the emergent communication protocols is commonly measured through topographic similarity (TOPSIM; Brighton and Kirby, 2006). It measures the topographic relation between meanings and signals, conceptually it gauges whether similar meanings map to similar signals. This metric was first used in recent computational simulations by Lazaridou et al. (2018) and has been used in a large body of work since. Yet, the interpretation of linguistic properties emerging in simulations remains challenging, since language protocols used among artificial agents often show critical mismatches with known properties of human languages (Galke et al., 2022; Lian et al., 2023) such as efficiency, word-order vs. case-marking biases, or compositional generalisation (see §2). Consequently, it is evident that their learning biases and signal-meaning mappings differ from those of humans. This underscores the critical need to obtain deeper insight into referential games in the language learning setting (Rita et al., 2022).

A possible explanation for these mismatches could stem from representational alignment, the degree of agreement between the internal representations of two information processing systems (Sucholutsky et al., 2023). To the best of our knowledge, representational alignment in emergent communication was first reported by (Bouchacourt and Baroni, 2018), who measured the degree to which

57

agents aligned their internal image interpretations (inter-agent alignment) by performing Representational Similarity Analysis (RSA; Kriegeskorte et al., 2008). Using RSA (§3), they showed that agents establish successful communication artificially by aligning their internal image representations while *losing* any relation to the images presented (image-agent alignment), enabling communication about noise input even though they were trained on real images. As such, their communication protocol captured not conceptual properties of the objects depicted in pictures, but most likely focused on non-human-like spurious image features (e.g., pixel intensities). While inter-agent alignment is not a problem per se, the loss of image-agent alignment is problematic for two reasons. First, for emergent communication simulations to provide meaningful insights into the emergence of natural human language, agent image representations must be grounded in the content of the images. Only then can we deduce *what* the agents communicate about and assess linguistic properties or their ability to generalise to novel concepts. Second, emergent communication setups have been proposed to fine-tune pre-trained (vision-)language models, aiming to enhance machine understanding of natural human language (Lazaridou and Baroni, 2020; Lowe et al., 2020; Steinert-Threlkeld et al., 2022; Zheng et al., 2024). In this context, maintaining substantial alignment between representations and images is crucial for preserving mutual understanding between machines and humans.
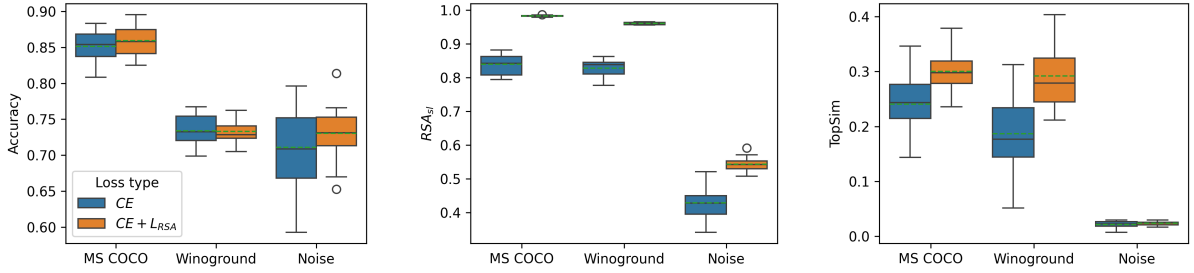
Representational alignment, however, did not receive the necessary attention since a host of papers appeared *after* the findings by Bouchacourt and Baroni in which results on referential games were reported without taking RSA into account (e.g. Lazaridou et al., 2018; Guo et al., 2019; Li and Bowling, 2019; Ren et al., 2020; Chaabouni et al., 2020; Dagan et al., 2021; Mu and Goodman, 2021; Chaabouni et al., 2022). Admittedly, some use attribute-value objects instead of real images as input. But importantly, in nearly all cases, neural agents must map inputs—whether attribute-value objects or image representations—onto agent-specific representations. Therefore the problem of inter-agent alignment *can always* occur and is *agnostic* to the input type. Although this warrants further analysis of earlier results, the field is already employing referential games in more complex simulations with real images (e.g. Dessì et al., 2021;

Chaabouni et al., 2022; Mahaut et al., 2024).

This work addresses the understudied alignment problem in standard referential game setups used in emergent communication. We train Reinforcement Learning (RL) agents equipped with a recent vision module (DinoV2; Oquab et al., 2024) to communicate about images. In addition to evaluating the agents on MS COCO (Lin et al., 2014) image pairs, we evaluate on noise pairs and image pairs sourced from the Winoground dataset (Thrush et al., 2022). The latter is explicitly created to gauge visio-linguistic compositional reasoning abilities of vision and language models. We first confirm that effective communication in the referential game relies on inter-agent alignment and then move on to our contributions. First, we find a strong correlation between the degree of inter-agent alignment and the TOPSIM metric. Our second contribution consists of a solution to the alignment problem by including an alignment penalty term to the loss, resulting in equivalent communicative success and higher TOPSIM whilst ensuring that the agents communicate about images instead of spurious features (Figure 1). We then argue to start evaluating emergent communication protocols on more strict tasks that directly target the intuition behind popular metrics to obtain a clearer understanding of the protocols. Overall, our results highlight the importance of representational alignment in simulations of language emergence and underscore the need to better understand the divergence in human and artificial language emergence.

## 2 Background

Most research in simulating emergent communication is modelled after the Lewis signalling game (Lewis, 1969) with a speaker and a listener agent. The speaker observes a state (e.g., an image) and sends a signal to the listener who acts based on this signal. In the case of the referential game, this means selecting a target among distractors. Both agents are rewarded for successful communication, meaning the listener points to the target object. The solution of this game requires the agents to have a shared protocol (i.e., an artificial language) which typically emerges when the agents learn based on trial and error over multiple games. This resembles how for humans, language learning and use impose constraints like pressures for learnability and compression that shape our language design (Kirby et al., 2014, 2015). Importantly, the emer-

(a) Communicative performance (Accuracy) on discriminating two images.

(b) Inter-agent representational alignment ($\mathrm{RSA}_{sl}$) between agent representations.

(c) Topographic similarity (TOPSIM) between the images and the messages.

Figure 1: Inference results for different datasets after training on MS COCO images. In (a) we see that agents can discriminate MS COCO images but struggle with discriminating Winoground images. In (b) we see the effect of the loss function on the degree of inter-agent representational alignment and (c) implies that according to the TOPSIM metric, messages are more structured if the alignment penalty is used. The presented results are across 15 seeds and use the best-performing parameters resulting from our parameter sweep, dashed green lines indicate averages.

gent language in this setup is also shaped by biases resulting from, for example, the agent architecture, loss function, and learning protocol (Rita et al., 2022). The current work uses the referential game: a variant of the Lewis signalling game extensively used to explore language evolution (e.g. Steels and Loetzsch, 2012; Kirby et al., 2015; Lazaridou et al., 2017; Kottur et al., 2017; Lazaridou et al., 2018; Kharitonov et al., 2020; Chaabouni et al., 2022).

An important challenge in emergent communication is that artificial learners often do not behave the same way as human learners in experimental settings. Some emergent protocols do not follow Zipf's law and thus are anti-efficient unless pressures for brevity are introduced (Chaabouni et al., 2019a), others do not show the word-order vs. case-marking trade-off found in human languages (Chaabouni et al., 2019b; Lian et al., 2021). Additionally, there is an ongoing debate on the degree to which the emergent languages allow for compositional generalisation (Lazaridou and Baroni, 2020; Conklin and Smith, 2023). It has been suggested to introduce communicative (e.g., alternating speaker/listener roles) and cognitive (e.g., memory) constraints (Galke et al., 2022) and use more natural settings to promote more human-like patterns of language emergence with neural agents (Kouwenhoven et al., 2022). Doing so changes the learning pressures to which the agents need to adapt and can recover initially absent linguistic phenomena of natural language in emergent languages (for a review see Galke and Raviv, 2024). An example of such work, investigating the word-order vs. case-marking trade-off, has succeeded in replicating this trade-off for neural learners (Lian

et al., 2023). Their setup differs from other work in that agents first learn a miniature language via supervised learning, and then optimise it for communicative success via RL, resulting in emergent languages that share linguistic universals with human language.

To enhance understanding of emergent communication in the Lewis game, Rita et al. (2022) decomposed the standard objective in Lewis games into two key components: a co-adaptation loss and an information loss. In doing so, they shed light on potential sources of overfitting and how they might hinder the emergence of structured communication protocols. They demonstrated that desired linguistic properties (e.g., compositionality and generalisability) emerge when they control the listener's ability to converge to the speaker agent (i.e., control for overfitting on the co-adaptation loss). While the co-adaptation loss has parallels to inter-agent alignment, their work does not address the alignment between the agents' image representation and the input features, which we deem crucial in developing grounded communication protocols.

Another challenge in emergent communication is the disentanglement of the underlying meanings of emergent languages. Earlier studies by Lazaridou et al. (2017) suggested that agents assign symbols to general conceptual properties of objects in images, rather than low-level visual features. However, as previously mentioned, follow-up work from Bouchacourt and Baroni (2018) showed this is not always the case. They found that agents align their agent-specific image representations without developing a language that captures conceptual properties depicted in the images. More-

over, agents lost any sense of meaningful within-category variation where two similar objects in human perception (e.g., two avocados) were observed as maximally dissimilar for the agents. In response to these findings, recent studies have implemented sanity checks, testing whether trained agents can communicate about noise (Dessi et al., 2021; Mahaut et al., 2024). However, to the best of our knowledge, there has been little attention to what we consider to be their main result: the alignment problem.

## 3 Representational alignment

Representational alignment is the degree of agreement between the internal representations of two information processing systems, whether biological or artificial. Even though widely recognised in cognitive science, neuroscience, and machine learning (Sucholutsky et al., 2023), representational alignment has not seen much interest in the field of emergent communication, except for the work by Bouchacourt and Baroni who analysed the referential game using RSA. This metric measures the alignment between two sets of numerical vectors, for example, image embeddings and agents' representations thereof. In practice, it is calculated by taking the pairwise (cosine) distances between vectors of a set and calculating the Spearman rank correlation between these distances.

In this paper, we also use RSA to operationalise representational alignment. Given the speaker image representations $r_s$ of the DinoV2 input embeddings $i$ and $r_l$ as the same images represented in the listener representation space, we compute the pairwise cosine similarity between the representations for the speaker $s_s$ and for the listener $s_l$ and calculate Spearman's $\rho$ between $s_s$ and $s_l$. As such, this measures the degree of inter-agent alignment (RSA$_{sl}$) between image representations $s_s$ and $s_l$, relative to their input. Additionally, we use it to measure image-agent alignment between the speaker/listener image representations and the DinoV2 embeddings (RSA$_{si}$ and RSA$_{li}$). Importantly, alignment is *agnostic to the type of input*, being either images or attribute-value objects and can always happen when inputs are projected onto agent-specific representations.

Intuitively, a high inter-agent RSA$_{sl}$ value can be interpreted as agents with *similar* representations for similar images. Importantly, this can have two causes: both agents' image representations

either *maintain* a relation to the image input, or *lose* this relation. While the former is desirable, the latter means that the agents are not communicating about the same high-level image features but are likely communicating about non-human-like spurious features. A low RSA$_{sl}$ value entails that the agents have developed *different* interpretations for the same image. While this may well be similar to the question of whether people have different perceptual experiences of colour (Locke, 1847), in the case of emergent communication, the agents should develop a grounded vocabulary with overlapping concept-level properties if we wish machines to have more natural understanding of human language. We use RSA 1) as a metric to reassess findings from Bouchacourt and Baroni and 2) as an auxiliary loss to mitigate the alignment problem and ensure that the agents communicate about image features.

## 4 Methods

The standard referential game is used as provided by the EGG framework (Kharitonov et al., 2021). Doing so ensures our findings are representative of the widely-used setup, rather than being influenced by specific design decisions. The game is implemented as a multi-agent cooperative RL problem where a speaker and a listener communicate to discriminate a target image from two shuffled distractor images. The speaker receives a target $t$ and generates a message $m$ of at most length $L$, using vocabulary $V$. Importantly, the messages and symbols have no a priori meaning but are assumed to obtain meaning and become grounded during the game. Once meaningful, the symbols are ideally combined in a structured manner to create compositional messages that express more complex meanings. Using message $m$, the listener guesses the target $\hat{t}$. Communicative success is defined as $\hat{t} = t$, meaning that the listener has correctly identified the target image among the candidate images.

### 4.1 Agents

Agents contain a language and a vision module. The latter consists of a frozen pre-trained visual network (DinoV2) and a learned agent-specific representation layer. While difficult to know what conceptual image features are present in DinoV2 embeddings, they have demonstrated capability in semantic segmentation tasks (Oquab et al., 2024), which is similar to the agents' objective. In contrast

to the hybrid structure of the vision module, the language module is entirely trained from scratch.

*The speaker* agent processes images by applying a linear transformation to the image embeddings, followed by batch normalisation, to create its agent-specific image representation $r_s$. Its language module embeds this representation and passes it through a single-layer Gated Recurrent Unit (GRU; Cho et al., 2014) that spells out messages to describe the target. *The listener* receives the message and the distractor images. It encodes the message into an embedding using another single-cell GRU layer. Additionally, a listener image representation $r_l$ is obtained for each image by applying a linear transformation and batch normalisation. Subsequently, temperature-weighted (temperature defaults to 0.1) cosine scores construct a multi-modal representation between the image and message representation (Dessi et al., 2021), where a higher probability should be assigned to the target image.

## 4.2 Optimisation

Communicative success ($\hat{t} = t$) is used to optimise the trainable parameters of both agents. The listener minimises cross-entropy ($ce$) loss using stochastic gradient descent, amounting to supervised learning. The $ce$ loss is calculated over the listeners' target distribution, thus providing direct pressure for communicative success. At inference, the candidate image with the highest probability is chosen as the target $\hat{t}$. The gradients required to optimise the speaker are calculated using the REINFORCE (Williams, 1992) update rule as each generated symbol must be assigned a loss. Following common practice (Rita et al., 2024), entropy regularisation (Mnih et al., 2016) is added to the loss to maintain exploration in message generation.

In addition to the conventional $ce$ loss, we introduce an alignment loss ($ce + \text{RSA}$) that includes an alignment penalty term to enforce high inter-agent and image-agent alignment. The term

$$L_{\text{RSA}} = (1 - \text{RSA}_{sl}) + (1 - \text{RSA}_{si}) + (1 - \text{RSA}_{li})$$

is added to the $ce$ loss with equal importance. We use torchsort (Blondel et al., 2020) to calculate $L_{\text{RSA}}$ such that the entire loss term is differentiable. Importantly, $L_{\text{RSA}}$ is not influenced by communicative success and does not interact with the $ce$ loss (Appendix C). Only adding $\text{RSA}_{sl}$ to the $ce$ loss is not sufficient as high inter-agent alignment can be achieved while *losing* image-agent alignment (see §3). We therefore also include $\text{RSA}_{si}$



Figure 2: Exemplar pairs of each dataset used for evaluation. Left: an image pair from MS COCO. Middle: A Winoground example. Right: A Gaussian noise pair. All images are cropped for display purposes.
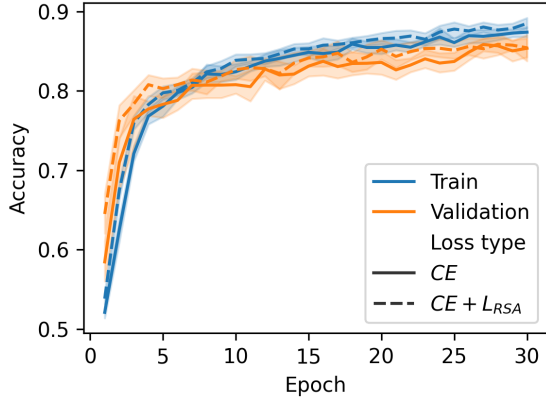
and $\text{RSA}_{li}$ to ensure that the agents communicate about the content displayed in the images. Including $\text{RSA}_{sl}$ entails that representational information is shared between the agents, thus differing from how humans interact. Yet, ranking the speaker and listener representations in calculating $\text{RSA}_{sl}$ bears *some* resemblance to projecting beliefs upon the interpretations of the other communicative partner. The current solution should be seen as a step towards more grounded vocabularies prone to refinements such as cognitive plausibility. We train for 30 epochs regardless of the loss used. The hyperparameters (Appendix B) that resulted in the best validation accuracy across 42 different communication channel capacities (Appendix A) were used for our findings.
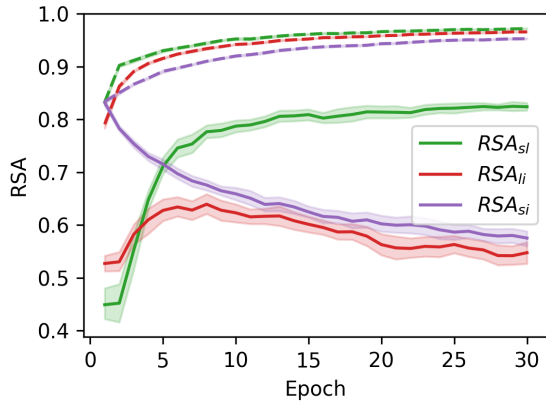
## 4.3 Data

Agents are trained to discriminate MS COCO images but tested on three different datasets (Figure 2) to assess out-of-distribution (o.o.d.) performance.

**MS COCO** – We use a subset of 1200 images from the MS COCO 2017 validation set to train and test the agents using an 80/20 split. To obtain this subset, we first select the categories that contain more than 100 images (12 categories) and subsequently sample 100 images for each supercategory present in the resulting set of images. The distractor images are sampled from the same category to ensure that there is *some* relevance to the target image. Importantly, sampling distractor images is done for each batch, meaning targets have different distractors at each epoch.

**Winoground** – The Winoground dataset (Thrush et al., 2022) was created to assess the visio-

(a) Learning curves for the MS COCO dataset on train and validation data.



(b) Representational alignment between agent image representations (green) and between the image and the sender/listener representations (purple, red).

Figure 3: In (a) we see that the agents learn to communicate successfully without overfitting on train data. In (b) we see that the alignment problem occurs with the *ce* but not the *ce* + RSA loss. Line style indicates the loss type. Data is averaged over 15 seeds, areas indicate the 95% confidence intervals.

linguistic compositional reasoning abilities of vision and language models. Here, we repurpose it as a proxy for the agents' ability to endow in compositional reasoning for image-based settings. The dataset contains 800 images and corresponding captions, comprising 400 Winoground pairs. Image-caption pairs were included when the captions share the same words but are of different *compositions*, implying completely different semantics (e.g., "a tree smashed into a car" versus "a car smashed into a tree" in Figure 2 (middle)). We only use the image pairs, not the captions. Crucially, this task differs from MS COCO since the image pairs are *fixed*, *conceptually similar* and meant to be discriminative if the agents' language allows

for compositional reasoning and is grounded in the visual modality.

**Noise** – Following Bouchacourt and Baroni (2018), we test whether agents can communicate about Gaussian noise ($\mu = 0, \sigma = 1$) pairs when trained on real images. Being able to do so would imply that messages communicate about spurious instead of high-level concept features.

### 4.4 Metrics

The performance of our agents is assessed by communicative success (accuracy) and RSA (§3) measures alignment. The degree of compositionality in the emergent language is assessed through the TOPSIM metric. Other metrics for compositionality like positional disentanglement, bag-of-symbols disentanglement (Chaabouni et al., 2020) are not straightforward due to the continuous nature of the image embeddings.

## 5 Results

### 5.1 Communicative success

Unsurprisingly, results show that agents can successfully disambiguate between image pairs from MS COCO using an emergent language (Figure 3a). Notably, we also confirm previous observations by (Bouchacourt and Baroni, 2018) that agents trained on real images can communicate about Gaussian noise (Figure 1a). Thus again suggesting that the messages convey spurious features rather than concept-level information. Interestingly, their performance on Gaussian noise is comparable to the performance on Winoground pairs, which requires the messages to capture concept-level properties. Thus revealing the difficulty of discriminating between strict pairs of conceptually similar images. The observed decrease in o.o.d. performance aligns with findings from other studies, such as Lazaridou et al. (2018) and Conklin and Smith (2023).

### 5.2 The alignment problem

The solid lines in Figure 3b clearly show that inter-agent alignment increases while alignment sensitivity to image features decreases for both agents. In principle, it is not a problem that the agents' image representations align. However, it is problematic when the alignment between the image embeddings and the image representations declines. Ablations across different channel capacities (§A) and pre-trained vision modules (§D) showed that

these trends appear consistently and are not influenced by the capacity or type of vision model. In addition to the communicative success on Gaussian noise, this re-confirms that the agents do not learn to extract concept-level information from the image embeddings but instead solve this task differently.
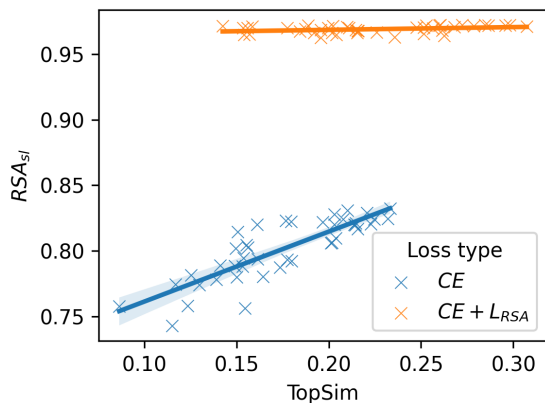
### 5.3 TOPSIM and representational alignment



Figure 4: The relationship between TOPSIM and inter-agent alignment ($\text{RSA}_{sl}$) for both loss types.

Earlier findings show mixed results on the relationship between TOPSIM and generalisation in image-based settings, TOPSIM was either related to generalisation (Chaabouni et al., 2022) or not (Rita et al., 2022). Our results indicate that generalisation and TOPSIM are correlated with both $ce$ ($r = .856$, $p < .001$) and $ce + \text{RSA}$ ($r = .767$, $p < .001$) losses. Meaning that more structured languages enable better communication on unseen validation pairs. Moreover, we find a strong positive relationship between $\text{RSA}_{sl}$ and TOPSIM ($r = .838$, $p < .001$) in the $ce$ (Figure 4). This relation is also present in the $ce + \text{RSA}$ setup ($r = .408$, $p = .001$), but is decoupled from TOPSIM given the (very) small spread ($\sigma = .003$) of $\text{RSA}_{sl}$. We do not observe an influence of inter-agent alignment on the number of uniquely produced messages.

### 5.4 Mitigating the alignment problem

We now focus on the $ce + \text{RSA}$ setup which was introduced to ensure that the agents maintain alignment with the image embeddings. Figure 3b shows that this is the case: inter-agent alignment *and* agent-image alignment increase during training and remain high at inference. However, there does not seem to be a benefit for communicative success at inference time (Figure 1). This is because the

alignment penalty only forces agents to represent images similarly to the image embeddings and is independent of the cross-entropy loss used to assess the success of communication (Appendix C). In the case of noise images, we still observe the above-chance performance, suggesting that communication between the agents still occurs in an artificial manner.

The alignment penalty also leads to increased TOPSIM, indicating a higher level of structure (Figure 1c) and strengthens our finding that TOPSIM and inter-agent alignment are related. Suggesting that the observed variations in TOPSIM, whether higher or lower, as noted in previous studies (e.g. Kottur et al., 2017; Chaabouni et al., 2020), should not be interpreted without considering alignment since they may be attributable to this underlying artefact rather than alterations to the original setup.

When tested on more strict Winoground pairs, communicative success does not improve as a result of using the alignment penalty (Figure 1a). Given the correlation between TOPSIM and generalisation, this is surprising since the higher degree of TOPSIM should imply that the language is more structured. Moreover, both, $\text{RSA}_{si}$ and $\text{RSA}_{li}$ have not drifted away from the image features. This combination, *in theory*, should be ideal for discriminating image pairs from the Winoground dataset since it was designed to be discriminative with compositional visio-linguistic reasoning. However, in *practice* this is not the case.

## 6 Discussion

In this work, we revisited the representational alignment problem in a common setup used in emergent communication and proposed a solution to this underrepresented problem. We corroborated earlier findings by showing that agents align their image representations and rely on spurious image features instead of human-like concept-level information (Bouchacourt and Baroni, 2018). We then showed that inter-agent alignment strongly correlates with the commonly used TOPSIM metric. Our solution to the alignment problem involves an alignment penalty that forces the agents to remain aligned with the input features and mitigates the alignment problem without decreasing communicative success. Finally, when agents are tested on more challenging Winoground pairs we observed reasonable but lower performance regardless of whether image representations were similar to the image em-

beddings or not. With this work, we hope that the alignment problem will receive more attention in the field of emergent communication, as is already the case in adjacent fields (Sucholutsky et al., 2023).

## 6.1 Importance of representational alignment

It is common practice in simulations of emergent communication to process (visual) inputs into an agent-specific hidden representation and update their weights simultaneously (e.g. Lazaridou et al., 2017; Bouchacourt and Baroni, 2018; Chaabouni et al., 2019a, 2020; Rita et al., 2022). As such, inter-agent alignment, *irrespective of the input form*, likely happens in other simulations too. This phenomenon is therefore potentially widespread and perhaps the cause for findings that are at odds with experimental findings. While it is not always the case that the representation structure we *expect* to help solve a task will do so (e.g. Montero et al., 2021; Xu et al., 2022), such discrepancies may hinder the use of emergent communication models in developing a more natural understanding of human languages and leave them less suitable for directly simulating language evolution phenomena. Especially if we want machine representations of natural language to align with human representations (Sucholutsky et al., 2023). RSA should therefore be used to rule out, or at the bare minimum report about, representational alignment in the future.

## 6.2 TOPSIM and representational alignment

Measuring representational alignment using RSA is similar to how TOPSIM measures the structure in messages. They differ in their inputs but both calculate the Spearman-ranked correlation between metric-agnostic pairwise distances. Crucially, the input makes all the difference, the inputs for RSA are from both agents and are trained independently, whilst TOPSIM only assesses the relation between the fixed inputs and learned output. Despite the similarities, the metrics thus describe different phenomena and are rarely reported simultaneously.

We hypothesise that the relationship between TOPSIM and inter-agent representational alignment is a by-product of the setup, which in essence implies that the listener has to align its representation $r_l$ to the speaker representation $r_s$ (Rita et al., 2022). It has to do so using only the speakers' messages, being an abstraction of $r_s$. A solution to this problem is to align representations, which eases the listeners' training objective. If the speaker

consistently produces structured messages during training, aligning $r_l$ to $r_s$ is easier, thereby causing higher inter-agent alignment. Essentially, this renders TOPSIM to be an *indirect* metric for the rate of alignment, for which $RSA_{sl}$ is a *direct* metric. In the context of learnability, the relationship between TOPSIM and inter-agent alignment and the fact that alignment always occurs can be seen as reasons for why languages with higher TOPSIM are easier to learn (Li and Bowling, 2019; Cheng et al., 2023). This underscores the need to report inter-agent representational alignment to avoid conclusions drawn about the effect of specific interventions on TOPSIM which may be attributable to inter-agent alignment.

## 6.3 Targeted o.o.d. evaluations

An important implication of our findings concerns the standard practice of reporting o.o.d. accuracy where the agents are tested on unseen input after training (e.g. Auersperger and Pecina, 2022; Conklin and Smith, 2023). This should inform about the agents' ability to generalise from one dataset (e.g., MS COCO) to another dataset (e.g., the Winoground pairs) much like human language allows us to talk about an infinite number of situations. Crucially, this overlooks the representational alignment problem in that we do not know *what* the agents are precisely generalising about. This problem can be mitigated with the alignment penalty to assess generalisation more directly or at least should be taken into consideration.

We assess o.o.d. performance on the more challenging Winoground pairs as a proxy for the agents' ability to endow in compositional reasoning for image-based settings. Good performance on the Winoground dataset requires a grounded language that can be used to create compositional messages since the objects and their underlying relations need to be described. In general, we suggest to start evaluating simulations of referential games on targeted strict tasks, like probing state-of-the-art vision language models on e.g., visio-compositional (Thrush et al., 2022; Diwan et al., 2022; Hsieh et al., 2023; Ray et al., 2023) or spatial (Kamath et al., 2023) reasoning. Re-purposing such datasets can reveal more directly whether agents develop the attested communicative abilities that are trivial to humans without having to rely on metrics. Our results illustrate this through a shortcoming of the TOPSIM metric. We observed that agents still struggle with distinguishing pairs of *conceptually*

*similar* Winoground images even though TOPSIM is higher with the alignment penalty. If the language protocol were to communicate concept-level information *and* compositional messages were created, we should not observe this struggle, meaning that the emerged protocols do not enable human-like communicative success.

Interestingly, the o.o.d. performance remains substantially above chance in the $ce + $ RSA setting. Given that MS COCO is not a dataset for learning to model compositionality, this delineates the limits of what can be achieved qua performance based on MS COCO image features in the Winoground context. Nevertheless, this leaves open the question of above-chance performance on Gaussian noise with the $ce + $ RSA loss. A tentative explanation is that the higher inter-agent alignment on noise input ($M_{ce} = .428$, $M_{ce+\text{RSA}_{sl}} = .543$, $t = -8.71$, $p < .001$) alleviates part of the problem (Figure 1b). To validate this, future experiments should involve controlling the prior distributions of the agents' image encoders by training their vision modules on different data. Doing so ensures that they have to communicate about novel objects and cannot rely on similar representations.

## 7 Conclusion

This paper revisits the underrepresented alignment problem present in the referential game often used in simulations of emergent communication. Specifically, we focused on the problem of increasing alignment between agent-image representations in combination with a decreasing alignment between the input and agent representations. We first confirmed that agents align their image representations while losing connection to their input, meaning that the emergent languages do not appear to encode human-like visual features. We then showed that, in the common setup, inter-agent alignment is related to topographic similarity, and argued that this renders TOPSIM an *indirect* metric of the rate of inter-agent alignment. To further investigate the effects of alignment, we introduced an alignment penalty to mitigate the alignment problem and showed that the communicative ability on a strict compositionality benchmark did not improve, leaving the question of inducing compositional generalisation in emergent communication for images unsolved. Our findings underscore the need to better understand the divergence between human and artificial language emergence within the prevalent

referential setup and highlight the importance and potential impact of representational alignment. We hope that future work rules out or at least reports about representational alignment.

## 8 Limitations

Our work has a few notable limitations. First, it only involves the referential game. Another popular variant, the reconstruction game (e.g. Chaabouni et al., 2019a, 2020; Lian et al., 2021; Conklin and Smith, 2023), requires the listener to reconstruct the input object based on the speakers' message. Since this setup has a different objective and presents different learning biases, it may have different results. We still expect the results to be similar as there is no pressure to retain alignment between the image input and agent representation. It would, however, be interesting to investigate whether the language protocol in this scenario is more structured than in the referential game.

Another limitation in our setup is that we only consider the scenario with two agents, which may be a requirement for alignment to be possible. Since experiments with human participants show that larger communities create more systematic languages (Raviv et al., 2019b), simulations on emergent multi-agent communication with populations of agents are also conducted, but these yield mixed results. The emergent communication protocols oftentimes do not evolve to be more structured unless explicit pressures such as population diversity or emulation mechanisms are introduced (Rita et al., 2022; Chaabouni et al., 2022). Michel et al. (2023) however, showed that population setups can result in more compositional languages if agent pairs are trained in a partitioned manner to prevent co-adaptation. Despite the mixed results, we believe that emergent communication with populations of agents is ecologically more valid and could result in different alignment effects. Much like how Tieleman et al. (2019) showed that autoencoders encode better concept category representations when they learn representations in a community-based setting with multiple encoders and decoders collectively.

The final limitation of our study regards its scale. While simulations of emergent communication are typically conducted on relatively small-scale datasets, human language emergence is accompanied by rich and diverse multi-modal experiences. Recent results in the field of computer vision suggest that dataset diversity and scale are

the primary drivers of alignment to human representations (Conwell et al., 2023; Muttenthaler et al., 2023). As such, this key difference between the setting of artificial emergent communication and human language emergence can drive the observed differences in representations. Due to the difficulty of interpreting these representations, we see this as another reason to evaluate emergent protocols on more strict datasets with clear pragmatic value for humans.

# References

Michal Auersperger and Pavel Pecina. 2022. Defending compositionality in emergent languages. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*, pages 285–291, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

Mathieu Blondel, Olivier Teboul, Quentin Berthet, and Josip Djolonga. 2020. Fast differentiable sorting and ranking. In *International Conference on Machine Learning*, pages 950–959. PMLR.

Diane Bouchacourt and Marco Baroni. 2018. How agents see things: On visual representations in an emergent language game. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 981–985, Brussels, Belgium. Association for Computational Linguistics.

Henry Brighton and Simon Kirby. 2006. Understanding linguistic evolution by visualizing the emergence of topographic mappings. *Artificial Life*, 12(2):229–242.

Rahma Chaabouni, Eugene Kharitonov, Diane Bouchacourt, Emmanuel Dupoux, and Marco Baroni. 2020. Compositionality and generalization in emergent languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4427–4442, Online. Association for Computational Linguistics.

Rahma Chaabouni, Eugene Kharitonov, Emmanuel Dupoux, and Marco Baroni. 2019a. Anti-efficient encoding in emergent communication. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Rahma Chaabouni, Eugene Kharitonov, Alessandro Lazaric, Emmanuel Dupoux, and Marco Baroni. 2019b. Word-order biases in deep-agent emergent communication. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5166–5175, Florence, Italy. Association for Computational Linguistics.

Rahma Chaabouni, Florian Strub, Florent Altché, Eugene Tarassov, Corentin Tallec, Elnaz Davoodi, Kory Wallace Mathewson, Olivier Tieleman, Angeliki Lazaridou, and Bilal Piot. 2022. Emergent communication at scale. In *International Conference on Learning Representations*.

Emily Cheng, Mathieu Rita, and Thierry Poibeau. 2023. On the correspondence between compositionality and imitation in emergent neural communication. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12432–12447, Toronto, Canada. Association for Computational Linguistics.

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.

Henry Conklin and Kenny Smith. 2023. Compositionality with variation reliably emerges in neural networks. In *The Eleventh International Conference on Learning Representations*.

Colin Conwell, Jacob S. Prince, Kendrick N. Kay, George A. Alvarez, and Talia Konkle. 2023. What can 1.8 billion regressions tell us about the pressures shaping high-level visual representation in brains and machines? *bioRxiv*.

Gautier Dagan, Dieuwke Hupkes, and Elia Bruni. 2021. Co-evolution of language and agents in referential games. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2993–3004, Online. Association for Computational Linguistics.

Bart de Boer. 2006. Computer modeling as a tool for understanding language evolution. In Nathalie Gontier, Jean Paul Van Bendegem, and Diederik Aerts, editors, *Evolutionary Epistemology, Language and Culture*, pages 381–406. Springer, Dordrecht.

Roberto Dessi, Eugene Kharitonov, and Baroni Marco. 2021. Interpretable agent communication from scratch (with a generic visual processor emerging on the side). In *Advances in Neural Information Processing Systems*, volume 34, pages 26937–26949. Curran Associates, Inc.

Anuj Diwan, Layne Berry, Eunsol Choi, David Harwath, and Kyle Mahowald. 2022. Why is winoground hard? investigating failures in visuolinguistic compositionality. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2236–2250, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Lukas Galke, Yoav Ram, and Limor Raviv. 2022. Emergent communication for understanding human language evolution: What's missing? In *Emergent Communication Workshop at ICLR 2022*.

Lukas Galke and Limor Raviv. 2024. Emergent communication and learning pressures in language models: a language evolution perspective. *arXiv preprint arXiv:2403.14427*.

Shangmin Guo, Yi Ren, Serhii Havrylov, Stella Frank, Ivan Titov, and Kenny Smith. 2019. The emergence of compositional languages for numeric concepts through iterated learning in neural agents.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.

Charles F Hockett. 1959. Animal" languages" and human language. *Human Biology*, 31(1):32–39.

Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. 2023. Sugarcrepe: Fixing hackable benchmarks for vision-language compositionality. In *Advances in Neural Information Processing Systems*, volume 36, pages 31096–31116. Curran Associates, Inc.

Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023. What's "up" with vision-language models? investigating their struggle with spatial reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9161–9175, Singapore. Association for Computational Linguistics.

Eugene Kharitonov, Rahma Chaabouni, Diane Bouchacourt, and Marco Baroni. 2020. Entropy minimization in emergent languages. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5220–5230. PMLR.

Eugene Kharitonov, Roberto Dessì, Rahma Chaabouni, Diane Bouchacourt, and Marco Baroni. 2021. EGG: a toolkit for research on Emergence of lanGuage in Games. https://github.com/facebookresearch/EGG.

Simon Kirby, Hannah Cornish, and Kenny Smith. 2008. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. *Proceedings of the National Academy of Sciences*, 105(31):10681–10686.

Simon Kirby, Tom Griffiths, and Kenny Smith. 2014. Iterated learning and the evolution of language.

Simon Kirby, Monica Tamariz, Hannah Cornish, and Kenny Smith. 2015. Compression and communication in the cultural evolution of linguistic structure. *Cognition*, 141:87–102.

Satwik Kottur, José Moura, Stefan Lee, and Dhruv Batra. 2017. Natural language does not emerge 'naturally' in multi-agent dialog. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2962–2967, Copenhagen, Denmark. Association for Computational Linguistics.

Tom Kouwenhoven, Tessa Verhoef, Roy de Kleijn, and Stephan Raaijmakers. 2022. Emerging Grounded Shared Vocabularies Between Human and Machine, Inspired by Human Language Evolution. *Frontiers in Artificial Intelligence*, 5.

Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. 2008. Representational similarity analysis-connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:4.

Angeliki Lazaridou and Marco Baroni. 2020. Emergent multi-agent communication in the deep learning era. *arXiv preprint arXiv:2006.02419*.

Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls, and Stephen Clark. 2018. Emergence of linguistic communication from referential games with symbolic and pixel input. In *International Conference on Learning Representations*.

Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. 2017. Multi-agent cooperation and the emergence of (natural) language. *5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings*, pages 1–11.

David Lewis. 1969. Convention: A philosophical study. *Cambridge, MA*.

Fushan Li and Michael Bowling. 2019. Ease-of-teaching and language structure from emergent communication. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Yuchen Lian, Arianna Bisazza, and Tessa Verhoef. 2021. The effect of efficient messaging and input variability on neural-agent iterated language learning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10121–10129, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Yuchen Lian, Arianna Bisazza, and Tessa Verhoef. 2023. Communication Drives the Emergence of Language Universals in Neural Agents: Evidence from the Word-order/Case-marking Trade-off. *Transactions of the Association for Computational Linguistics*, 11:1033–1047.

Tsung-Yi Lin, Michael Maire, Serge 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 13*, pages 740–755. Springer.

John Locke. 1847. *An essay concerning human understanding*, volume 114. Kay & Troutman.

Ryan Lowe, Abhinav Gupta, Jakob Foerster, Douwe Kiela, and Joelle Pineau. 2020. On the interaction between supervision and self-play in emergent communication. In *International Conference on Learning Representations*.

Matéo Mahaut, Francesca Franzon, Roberto Dessì, and Marco Baroni. 2024. Referential communication in heterogeneous communities of pre-trained visual deep networks.

Paul Michel, Mathieu Rita, Kory Wallace Mathewson, Olivier Tieleman, and Angeliki Lazaridou. 2023. Revisiting populations in multi-agent communication. In *The Eleventh International Conference on Learning Representations*.

Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1928–1937, New York, New York, USA. PMLR.

Milton Llera Montero, Casimir JH Ludwig, Rui Ponte Costa, Gaurav Malhotra, and Jeffrey Bowers. 2021. The role of disentanglement in generalisation. In *International Conference on Learning Representations*.

Jesse Mu and Noah Goodman. 2021. Emergent communication of generalizations. In *Advances in Neural Information Processing Systems*, volume 34, pages 17994–18007. Curran Associates, Inc.

Lukas Muttenthaler, Jonas Dippel, Lorenz Linhardt, Robert A. Vandermeulen, and Simon Kornblith. 2023. Human alignment of neural network representations. In *The Eleventh International Conference on Learning Representations*.

Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. 2024. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Limor Raviv, Antje Meyer, and Shiri Lev-Ari. 2019a. Compositional structure can emerge without generational transmission. *Cognition*, 182:151–164.

Limor Raviv, Antje Meyer, and Shiri Lev-Ari. 2019b. Larger communities create more systematic languages. *Proceedings of the Royal Society B*, 286(1907):20191262.

Arijit Ray, Filip Radenovic, Abhimanyu Dubey, Bryan Plummer, Ranjay Krishna, and Kate Saenko. 2023. Cola: A benchmark for compositional text-to-image retrieval. In *Advances in Neural Information Processing Systems*, volume 36, pages 46433–46445. Curran Associates, Inc.

Yi Ren, Shangmin Guo, Matthieu Labeau, Shay B. Cohen, and Simon Kirby. 2020. Compositional languages emerge in a neural iterated learning model. In *International Conference on Learning Representations*.

Mathieu Rita, Paul Michel, Rahma Chaabouni, Olivier Pietquin, Emmanuel Dupoux, and Florian Strub. 2024. Language evolution with deep learning. *arXiv preprint arXiv:2403.11958*.

Mathieu Rita, Corentin Tallec, Paul Michel, Jean-Bastien Grill, Olivier Pietquin, Emmanuel Dupoux, and Florian Strub. 2022. Emergent communication: Generalization and overfitting in lewis games. In *Advances in Neural Information Processing Systems*, volume 35, pages 1389–1404. Curran Associates, Inc.

Reinhard Selten and Massimo Warglien. 2007. The emergence of simple languages in an experimental coordination game. *Proceedings of the National Academy of Sciences*, 104(18):7361–7366.

Kenny Smith. 2022. How Language Learning and Language Use Create Linguistic Structure. *Current Directions in Psychological Science*, 31(2):177–186.

Luc Steels and Martin Loetzsch. 2012. The Grounded Naming Game. *Experiments in cultural language evolution*, 3:41–59.

Shane Steinert-Threlkeld, Xuhui Zhou, Zeyu Liu, and C. M. Downey. 2022. Emergent communication fine-tuning (EC-FT) for pretrained language models. In *Emergent Communication Workshop at ICLR 2022*.

Ilia Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim, Bradley C. Love, Erin Grant, Iris Groen, Jascha Achterberg, Joshua B. Tenenbaum, Katherine M. Collins, Katherine L. Hermann, Kerem Oktar, Klaus Greff, Martin N. Hebart, Nori Jacoby, Qiuyi Zhang, Raja Marjieh, Robert Geirhos, Sherol Chen, Simon Kornblith, Sunayana Rane, Talia Konkle, Thomas P. O'Connell, Thomas Unterthiner, Andrew K. Lampinen, Klaus-Robert Müller, Mariya Toneva, and Thomas L. Griffiths. 2023. Getting aligned on representational alignment.

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 (CVPR)*, pages 5238–5248.

Olivier Tieleman, Angeliki Lazaridou, Shibl Mourad, Charles Blundell, and Doina Precup. 2019. Shaping representations through communication: community size effect in artificial learning systems. *arXiv preprint arXiv:1912.06208*.

Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256.

Zhenlin Xu, Marc Niethammer, and Colin A Raffel. 2022. Compositional generalization in unsupervised compositional representation learning: A study on disentanglement and emergent language. In *Advances in Neural Information Processing Systems*, volume 35, pages 25074–25087. Curran Associates, Inc.

Chenhao Zheng, Jieyu Zhang, Aniruddha Kembhavi, and Ranjay Krishna. 2024. Iterated learning improves compositionality in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13785–13795.

## A Channel capacity

To test to what degree communicative success, TOPSIM, and representational alignment are confounded with the communication channel capacity, we ran simulations altering the vocabulary size ($V = \{3, 5, 10, 20, 40, 50, 100\}$) and message length ($L = \{2, 3, 5, 10, 50, 100\}$) resulting in 42 parameter settings per loss type.

Overall, performance is relatively independent of the chosen configuration, but vocabulary size influences success more than message length (Figure 5). The hyperparameters that resulted in the best validation accuracy (i.e., generalisation; Chaabouni et al., 2022) for the standard $ce$ setup were $V = 40$ and $L = 2$. These parameters are used to produce the results in the main paper. Contra expectations, the vocabulary size also influenced TOPSIM more than message length. It, especially in the case of $ce + L_{\text{RSA}}$, is higher when messages are shorter but have access to a larger vocabulary (Figure 6).

Figure 7 shows that, regardless of capacity, inter-agent alignment ($\text{RSA}_{sl}$) increases while image-agent alignment ($\text{RSA}_{si}$ and $\text{RSA}_{li}$) decreases with the $ce$ loss. Interestingly, $\text{RSA}_{sl}$ is agnostic to capacity but a larger vocabulary size, not message length, reduces the degree of drifting away from the input. We hypothesise this to result from lower pressure to compress rich continuous embeddings into smaller discrete vocabulary embeddings.
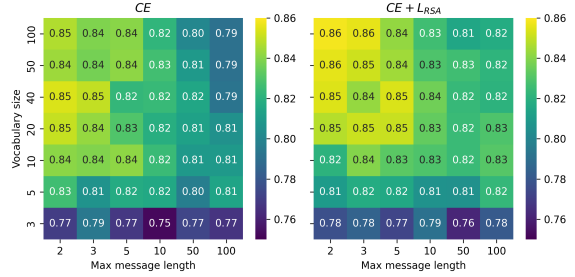


Figure 5: The validation accuracy as a dependent factor of the vocabulary size and maximum message length. Values are averages across 15 seeds.
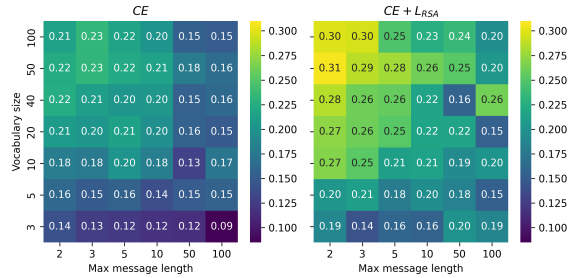


Figure 6: TOPSIM as a dependent factor of the vocabulary size and maximum message length. Values are averages across 15 seeds.

## B Best hyperparameters

The parameters used to run the experiments in the main paper were the following:

| Parameter | Value |
|---|---|
| Batch size | 32 |
| Optimiser | Adam |
| Learning Rate (S & L) | 0.01 & 0.001 |
| Vocabulary size ($V$) | 40 |
| Message length ($L$) | 2 |
| Hidden size (S & L) | 768 & 768 |
| Embedding size | 50 |
| Listener cosine temperature | 0.1 |
| Seeds | 16,22,41,56,67, 77,14,78,99,23, 82,40,51,37,62 |

Table 1: Best-performing parameters resulting from the parameter sweep.

## C Interaction of the alignment term on the cross-entropy loss

To ensure that there is no impact of the alignment penalty on the pressure for communicative success, we ablated the $L_{\text{RSA}}$ term of our proposed
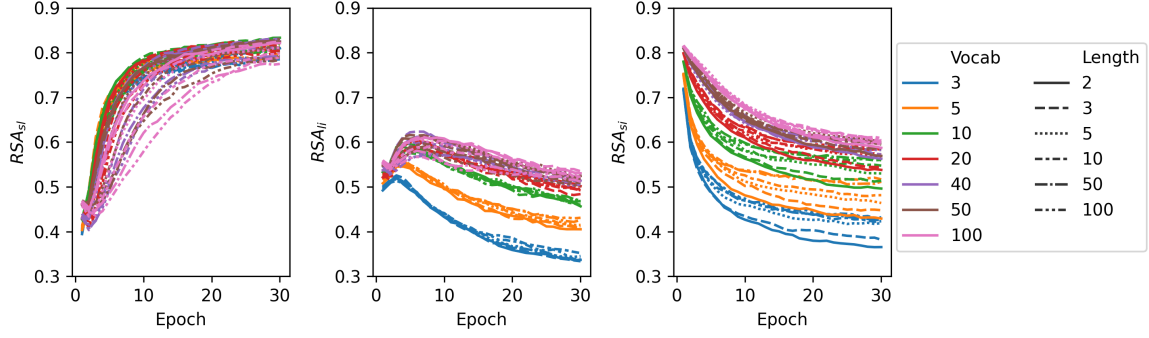
Figure 7: Representational alignment metrics averaged over 15 simulations with the standard *ce* loss. Regardless of channel capacity, representational alignment always occurs while losing relation to the input.

loss function and found that both, communicative success and *ce* are not affected by the alignment penalty (Figure 8). Corroborating that only the *ce* term provides pressure for successful communication (§5.4).
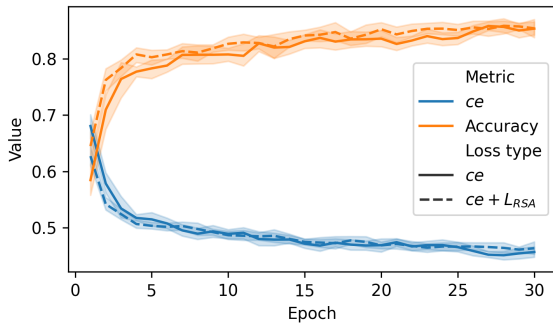


Figure 8: Learning curves (accuracy) and cross-entropy loss (*ce*) for both loss settings. There is virtually no effect of the auxiliary term $L_{\mathrm{RSA}}$ on the cross entropy loss or communicative success.

## D  Pre-trained vision modules

Although it is in principle possible to train the vision module of the agents from scratch (Dessi et al., 2021), in our work, agents' perception stems from a pre-trained vision-language model. Although there is reason to believe that DinoV2 embeddings capture high-level, conceptual image features useful for discriminating image pairs (Oquab et al., 2024), we assessed the degree to which the alignment problem occurs for different pre-trained models despite encoding the same objects. We ran additional simulations using image features obtained from ResNet (He et al., 2016) and CLIP (Radford et al., 2021) for 6 different parameter settings with the *ce* loss function. Here we used the parameters that resulted in the best, worst, mean, and quantile validation

performance from the parameter sweep in appendix A (see Table 2), and a sensible setup with $V = 10$ and $L = 5$.

| Msg. Length ($L$) | Vocab. Size ($V$) | Vision |
|:---:|:---:|:---:|
| 2 | 40 | |
| 3 | 10 | |
| 5 | 5 | DinoV2 |
| 5 | 10 | CLIP |
| 10 | 3 | ResNet |
| 50 | 100 | |

Table 2: The parameters for running additional simulations with CLIP and ResNet to assess the robustness of our results. Each combination was run for 15 different seeds. Note: results for the DinoV2 simulations are from the sweep.
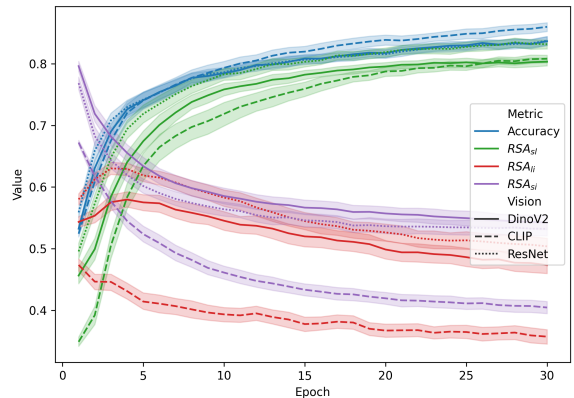


Figure 9: Learning curves (accuracy) and RSA metrics for different vision models averaged over 6 parameter settings with 15 seeds each. The representational alignment problem always occurs. Line style corresponds to the vision module used to obtain image embeddings and colour indicates the metric. Areas indicate the 95% confidence intervals.

Figure 9 shows clearly that inter-agent alignment

*increases* while agent-image alignment *decreases* for all models. In addition to the similar results reported by Bouchacourt and Baroni (2018) for VGG ConvNet embeddings, both 4096 and 1000 layers, we can confirm that the problem is agnostic to the input embeddings. Interestingly, agent representations drift most for CLIP embeddings. Nevertheless, the agents still develop a successful communication strategy, indicating that out-of-the-box CLIP embeddings are the least useful for agents in finding a (non-grounded) solution. No such differences are seen when the agents are trained with the additional alignment penalty term, inter-agent and image-agent alignment remain high for all models.