

What Makes Two Language Models Think Alike?

Anonymous ACL submission

Abstract

Do architectural differences significantly affect the way models represent and process language? We propose a new approach, based on metric-learning encoding models (MLEMs), as a first step to answer this question. The approach provides a feature-based comparison of how any two layers of any two models represent linguistic information. We apply the method to BERT, GPT-2 and Mamba. Unlike previous methods, MLEMs offer a *transparent* comparison, by identifying the specific linguistic features responsible for similarities and differences. More generally, the method uses formal, symbolic descriptions of a domain, and use these to compare neural representations. As such, the approach can straightforwardly be extended to other domains, such as speech and vision, and to other neural systems, including human brains.

1 Introduction

Marr’s hierarchy proposes a structured approach for describing information-processing systems using three levels (Figure 1; Marr, 2010): (1) computational, (2) algorithmic, and (3) implementational. The computational level defines the problem and the system’s goals. For example, a goal of a system could be to compute the sum of two numbers. The algorithmic level addresses the strategies used to solve the problem, detailing the step-by-step processes involved. For instance, one algorithm could involve digit-by-digit addition starting from the least significant digit, while another could involve repeated counting. There is therefore a one-to-many relationship between the computational and algorithmic levels (plain arrows) Finally, the implementational level concerns the physical realization of the system, such as how algorithms are executed within the brain’s neural architecture or a computer’s hardware. Similarly, there’s a one-to-many relation between the algorithmic and implementational levels (dashed arrows).

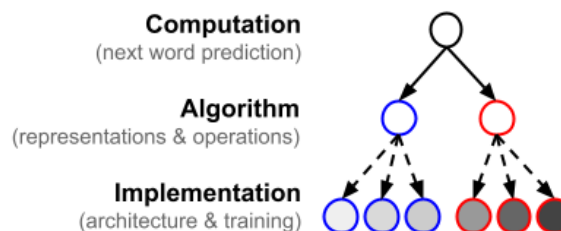


Figure 1: **Marr’s levels of analysis.** While language models may share the same computational goal (top level, next-word prediction), their architectures could differ substantially (bottom level). They therefore may or may not develop the same representations and algorithms (middle level) to perform the task.

Language models can be described along Marr’s three levels. At the computational level, most language models are trained on a next-word prediction task (but see, e.g., Gloeckle et al., 2024, for multi-token prediction). At the implementational level, different architectures (RNNs, Transformers, SSMs, etc.) can be implemented differently onto hardware. These architectural differences might lead to variations at the algorithmic level, despite sharing the same computational problem, or, conversely, they might converge on a similar algorithmic solution. In this work we ask whether language models with different architectures represent and process language in the same way?

To quantify this, one can start by computing second-order isomorphism between model representations of words (Shepard and Chipman, 1970), or use methods such as Canonical Correlation Analysis (CCA, Wu et al., 2020). Similarity can thus be computed for any pair of layers from any pair of models. However, similarity measures do not provide an *explanation* for what makes the representations and processing of text appear similar or dissimilar across models and layers.

To address this, we propose a novel approach based on metric-learning encoding models

068	(MLEMs, Jalouzet et al., 2024), which explains	118
069	model similarity by identifying the linguistic fea-	119
070	tures that underlie it. We illustrate the approach	120
071	using three different types of neural architectures:	121
072	encoder-based Transformer, decoder-based Trans-	122
073	former, and Mamba, quantifying their similarity	123
074	and providing a linguistic-feature-based compari-	124
075	son for each pair of model layers.	125
076	Overall, metric-learning encoding models use ex-	126
077	isting theoretical descriptions as a grid of analysis	127
078	of models and model comparisons. This approach	128
079	can thus naturally be extended beyond text to do-	
080	main domains such as speech and vision. They can then	
081	use any symbolic theory there to compare any two	
082	neural models, including artificial neural models	
083	(different architectures or different instantiations	
084	of the same one) as well as human and non-human	
085	animal brains.	
086	2 Related Literature	
087	In previous work to quantify similarity between	
088	representations of two neural systems, a central	
089	approach is based on <i>second-order isomorphism</i>	
090	(Shepard and Chipman, 1970). Second-order iso-	
091	morphism suggests that while the representations	
092	of two systems belong to different spaces, the simi-	
093	larity between them can be quantified by compar-	
094	ing the pairwise distances within each neural	
095	space, thus ‘second-order’ similarity. Second-order	
096	isomorphism has been used to compare represen-	
097	tations of two artificial neural networks (Laakso	
098	and Cottrell, 2000 ; Mehrer et al., 2020), or of	
099	two brains, where it is also known as Representa-	
100	tional Similarity Analysis (RSA; Kriegeskorte	
101	et al., 2008 ; Abnar et al., 2019).	
102	Several other similarity measures between repre-	
103	sentations of different models have been proposed	
104	in previous work, including linear regression (Adri-	
105	ana et al., 2015), canonical correlation analysis	
106	(CCA; Raghu et al., 2017 ; Morcos et al., 2018 ; Wu	
107	et al., 2020 ; Belinkov and Glass, 2018), statistical	
108	shape analysis (Williams et al., 2024), functional	
109	behaviors on downstream tasks (e.g., Alain and	
110	Bengio, 2018), and Dynamic Similarity Analysis	
111	(DSA, Ostrow et al., 2023). However, such mea-	
112	sures do not directly provide an explanation for	
113	<i>why</i> two neural systems converge or differ in the	
114	way they represent information.	
115	Recently, metric-learning encoding models	
116	(MLEMs; Jalouzet et al., 2024) have been pro-	
117	posed as a method to examine the types of infor-	
	mation that predict neural distances between repre-	118
	sentations within a single neural system. MLEMs	119
	have shown their ability to identify which linguis-	120
	tic features most strongly predict neural distances	121
	in various layers of a model. Here, we leverage	122
	MLEMs to study the similarity between represen-	123
	tations in two different language models. This ap-	124
	proach offers a feature-based comparison of how	125
	two models represent linguistic information, and	126
	thereby explains the underlying factors driving the	127
	similarities and differences.	128
	3 General Setup	129
	Language Models We investigated the simi-	130
	larities between three different types of models:	131
	(1) GPT-2 (Radford et al., 2019), a decoder-based	132
	Transformer , (2) BERT (Devlin et al., 2018), an	133
	encoder-based Transformer , and (3) Mamba, an	134
	architecture based on a state-space model (Gu and	135
	Dao, 2023). We collected representations from	136
	each layer of the models for every word in our con-	137
	trolled dataset (see below). For words that are split	138
	into multiple tokens, we used the representation of	139
	the final token.	140
	Probing Data To study the neural encoding of	141
	linguistic features, we created a dataset, which con-	142
	tains a list of sentences and their corresponding list	143
	of linguistic features. Sentences and features were	144
	generated using a custom grammar to cover central	145
	linguistic features, such as grammatical number,	146
	gender or tense (Table S1).	147
	Metric-Learning Encoding Models (MLEMs)	148
	Metric-Learning Encoding Models (Jalouzet et al.,	149
	2024) start from the assumption that to effectively	150
	capture multivariate, distributed neural encoding of	151
	linguistic information, one should model <i>distances</i>	152
	among neural representations rather than individual	153
	activations of single units. Given a set of inputs	154
	(e.g., words), where each is represented along a set	155
	of features (e.g., tense, gender), the goal is to learn	156
	a <i>metric function</i> (aka, a distance function), which	157
	is defined over pairs of inputs and computed based	158
	on the features of the inputs only. The optimal such	159
	metric function is the one that minimizes the dif-	160
	ferences between the modeled distances among the	161
	inputs and the empirical (neural) ones (Figure S6).	162
	This optimal metric can be derived using standard	163
	metric-learning methods (Kulis et al., 2013).	164

4 Results

Feature-Importance Profiles To quantify similarity among models, we first ask which linguistic features best explain neural distances in each layer of a language model. For this, we computed *feature importance (FI)* based on Metric-Learning Encoding Models. That is, for each layer of a given language model, we computed which linguistic features (tense, grammatical number, etc.) predict neural distances among representations of words in the dataset. Specifically, we computed FI as the average decrease in Spearman correlation score of the trained MLEM on a left-out dataset when permuting a feature. We highlight several main observations in the results (Figure 3): first, part-of-speech is the dominant linguistic feature across layers of Transformer-based models. However, for Mamba, it is so only for the first and last layer. In Mamba, we observe a significant increase in the importance of word position at around layer 10 of the model. Finally, we note that the importance of the grammatical number feature tends to decrease from early to later layers in all models.

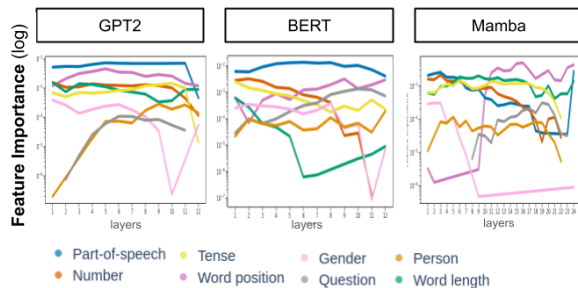


Figure 3: **Feature Importance Profiles.** The relative importance of linguistic features varies across layers and models

Feature-Based Similarity among Language Models We next asked, which language models most resemble each other in the way they represent linguistic information? We compared two approaches: *feature-based* and *feature-agnostic* similarity measures. For the former, we computed feature-based similarity based on the FI profiles from the MLEMs. Specifically, for each pair of layers, we computed the Kendall correlation coefficient, which quantifies to what extent the same linguistic features are dominant in both layers. Since linguistic feature with low importance (e.g., near zero) are not predictive of neural distances, it is desired that they will have a small effect on the similarity measure for two models. We therefore quantified similarity with a *weighted* version of Kendall correlation, which weighs feature importance based on their rank (Vigna, 2015). For comparison, for the feature-agnostic similarity measure, we followed a standard RSA approach (Kriegeskorte et al., 2008). Specifically, for each layer of a model, we first computed a dissimilarity matrix (DSM) for all words in the dataset. That is, for each layer, we computed the Euclidean distance among all pairs of stimuli presented to the model. Then, given two DSMs of two different layers, we computed the Spearman correlation between the upper triangles of the DSMs.

Figure 2 shows the resulting feature-based (Panel A) and the feature-agnostic (Panel B) matrices. To further visualize the results, the corresponding plots show each model layer in a shared 2D space, optimally preserving layer-wise similarity using Multi-Dimensional Scaling (MDS, Kruskal, 1964) analyses. Overall, the feature-based and feature-agnostic approaches agree on model simi-

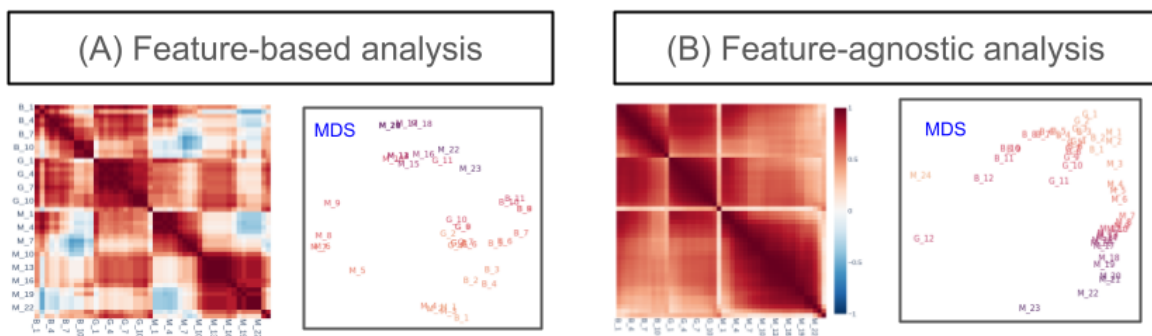


Figure 2: **Model Similarity.** (A) *Feature-based* similarity matrix corresponding to the pairwise correlations between feature-importance values. (B) *Feature-agnostic* similarity matrix based on raw Euclidean distances between word embeddings. The Multi-Dimensional Scaling representations of these distances are represented for both types of analyses (B stands for BERT, G for GPT2, and M for Mamba).

224 larity ($\rho_{Spearman} = 0.69$). However, feature-based
 225 similarity highlights specific differences between
 226 and across models. For example, for Mamba, a
 227 block structure appears, separating low and high
 228 layers of the model. The FI profiles for Mamba ex-
 229 plains this difference, given the sudden increase in
 230 FI of word position at around layer 10 of the model.
 231 This increase in FI is apparent when visualizing
 232 all word representations in the model for different
 233 layers. Indeed, word position strongly separates
 234 word representations at higher but not lower lay-
 235 ers (Figure 3). This illustrates the importance of a
 236 feature-based compared to feature-agnostic theory
 237 in explaining similarity, as we further investigate
 238 next.

239 **What Makes Two Language Models Think**
 240 **Alike?** With the feature-based similarity ap-
 241 proach, we are able to answer the question - why
 242 two model layers are similar or different in the
 243 way they represent linguistic information? Figure 4
 244 illustrates this by contrasting feature impor-
 245 tance of different pairs of layers - one with high
 246 ($\tau_{weighted} = 0.77$) and the other with low simi-
 247 larity ($\tau_{weighted} = -0.24$). These examples were
 248 chosen based on the minimal and maximal values
 249 of the similarity matrix (Figure 2A). For the case of
 250 high similarity (between GPT2 layer 8 and BERT
 251 layer 4), FI values of the two layers largely agree,
 252 lying on the diagonal of the scatter. In particu-
 253 lar, the most dominant feature in both layers is the
 254 same - part-of-speech (PoS). In accordance, the
 255 corresponding MDS plots (to the left and top of
 256 the scatter), show that word representations (color
 257 dots) are, indeed, well separated with respect to
 258 part-of-speech (see legend). In contrast, for the
 259 case of low similarity (between Mamba-layer-8
 260 and BERT-layer-4), FI values are mostly off di-
 261 agonal, in particular the most dominant one for PoS.
 262 In accordance, the corresponding MDS plot for
 263 Mamba does not separate word representations as
 264 well as in the two other cases.

265 **5 Summary and Conclusions**

266 While language models share the same computa-
 267 tional task (next-word prediction), architectural dif-
 268 ference might lead to differences in how different
 269 models represent and process language. Here, we
 270 presented a new approach to quantify such differ-
 271 ences. We illustrated the approach with three types
 272 of architecture, showing its utility in quantifying
 273 model similarity and, importantly, explaining it.

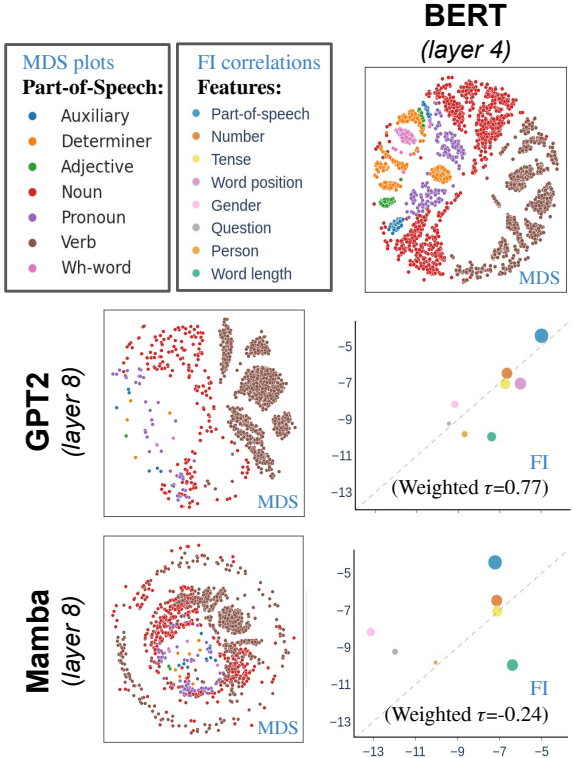


Figure 4: **Illustrating** how model/layers represent linguistic features. MDS plots of the representations, and pairwise comparison of the Feature Importance profiles.

274 For all pairs of model layers, we identified which
 275 linguistic features dominate word representations
 276 and whether they are the same or not across mod-
 277 els, as illustrated for the case of part-of-speech.
 278 Together, this shows the utility of feature-based
 279 approaches to study model similarity, providing
 280 theory-based explanations for why two models con-
 281 verge or diverge in the way they process natural
 282 language text. This approach could naturally be ex-
 283 tended to other domains, such as speech and vision.
 284 And it can be applied to compare neural systems,
 285 including artificial neural networks as we did here
 286 as well as human brains.

287 **Limitations**

288 For simplicity, when computing feature importance,
 289 we assumed that there are no interactions among
 290 linguistic features in predicting neural distances
 291 among sentence representations (i.e., assuming a
 292 diagonal weight matrix). However, such interac-
 293 tions are common in many problems, includ-
 294 ing in language. The framework of MLEMs al-
 295 lows a straightforward way to introduce interac-
 296 tions, while, in contrast to other approaches (such
 297 as RSA), it preserves the metric property of the

298	learned distances. Also, in MLEMs, we have only	Aarre Laakso and Garrison Cottrell. 2000. Content	350
299	included features that we consider essential to the	and cluster analysis: assessing representational simi-	351
300	list of words in the dataset (tense, grammatical	larity in neural systems. <i>Philosophical psychology</i> ,	352
301	number, etc.), as they were created by contrasting	13(1):47–76.	353
302	these dimensions. Future work can explore more	David Marr. 2010. <i>Vision: A computational investiga-</i>	354
303	exhaustive lists of features to describe and contrast	<i>tion into the human representation and processing of</i>	355
304	words, as well as larger datasets and the introduc-	<i>visual information</i> . MIT press.	356
305	tion of possible interactions among features.	Johannes Mehrer, Courtney J Spoerer, Nikolaus	357
306	References	Kriegeskorte, and Tim C Kietzmann. 2020. Individ-	358
307	Samira Abnar, Lisa Beinborn, Rochelle Choenni, and	ual differences among deep neural network models.	359
308	Jelle Zuidema. 2019. Blackbox meets blackbox: Rep-	<i>Nature communications</i> , 11(1):5725.	360
309	resentational similarity and stability analysis of neu-	Ari S. Morcos, Maithra Raghu, and Samy Bengio.	361
310	ral language models and brains. In <i>Proceedings of</i>	2018. Insights on representational similarity in neu-	362
311	<i>the ACL-Workshop on Analyzing and Interpreting</i>	ral networks with canonical correlation. In <i>Proceed-</i>	363
312	<i>Neural Networks for NLP</i> , pages 191–203.	<i>ings of the 32nd International Conference on Neu-</i>	364
313	Romero Adriana, Ballas Nicolas, K Samira Ebrahimi,	<i>ral Information Processing Systems, NIPS’18</i> , page	365
314	Chassang Antoine, Gatta Carlo, and Bengio Yoshua.	5732–5741, Red Hook, NY, USA. Curran Associates	366
315	2015. <i>Fitnets: Hints for thin deep nets</i> . <i>Proc. ICLR</i> ,	Inc.	367
316	2(3):1.	Mitchell Ostrow, Adam Eisen, Leo Kozachkov, and	368
317	Guillaume Alain and Yoshua Bengio. 2018. <i>Under-</i>	Ila Fiete. 2023. Beyond geometry: Comparing the	369
318	temporal structure of computation in neural circuits	370	
319	<i>probes</i> . <i>Preprint</i> , arXiv:1610.01644.	with dynamical similarity analysis . In <i>Advances in</i>	371
320	Yonatan Belinkov and James R. Glass. 2018. Analysis	<i>Neural Information Processing Systems</i> , volume 36,	372
321	methods in neural language processing: A survey .	pages 33824–33837. Curran Associates, Inc.	373
322	<i>Transactions of the Association for Computational</i>	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	374
323	<i>Linguistics</i> , 7:49–72.	Dario Amodei, Ilya Sutskever, et al. 2019. Language	375
324	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	models are unsupervised multitask learners. <i>OpenAI</i>	376
325	Kristina Toutanova. 2018. Bert: Pre-training of deep	<i>blog</i> , 1(8):9.	377
326	bidirectional transformers for language understand-	Maithra Raghu, Justin Gilmer, Jason Yosinski, and	378
327	ing. <i>arXiv preprint arXiv:1810.04805</i> .	Jascha Sohl-Dickstein. 2017. Svcca: Singular vector	379
328	Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière,	canonical correlation analysis for deep learning dy-	380
329	David Lopez-Paz, and Gabriel Synnaeve. 2024. Bet-	namics and interpretability . In <i>Advances in Neural</i>	381
330	ter & faster large language models via multi-token	<i>Information Processing Systems</i> , volume 30. Curran	382
331	prediction. <i>arXiv preprint arXiv:2404.19737</i> .	Associates, Inc.	383
332	Albert Gu and Tri Dao. 2023. Mamba: Linear-time	Roger N Shepard and Susan Chipman. 1970. Second-	384
333	sequence modeling with selective state spaces. <i>arXiv</i>	order isomorphism of internal representations:	385
334	<i>preprint arXiv:2312.00752</i> .	Shapes of states . <i>Cognitive Psychology</i> , 1(1):1–17.	386
335	Louis Jalouzet, Robin Sobczyk, Bastien Lhopitallier,	Sebastiano Vigna. 2015. A weighted correlation index	387
336	Jeanne Salle, Nur Lan, Emmanuel Chemla, and Yair	for rankings with ties. In <i>Proceedings of the 24th</i>	388
337	Lakretz. 2024. Metric-learning encoding models	<i>international conference on World Wide Web</i> , pages	389
338	identify processing profiles of linguistic features in	1166–1176.	390
339	bert’s representations . <i>Preprint</i> , arXiv:2402.11608.	Alex H. Williams, Erin Kunz, Simon Kornblith, and	391
340	Nikolaus Kriegeskorte, Marieke Mur, and Peter A. Ban-	Scott W. Linderman. 2024. Generalized shape met-	392
341	dettini. 2008. Representational similarity analysis	rics on neural representations. In <i>Proceedings of the</i>	393
342	– connecting the branches of systems neuroscience .	<i>35th International Conference on Neural Information</i>	394
343	<i>Frontiers in Systems Neuroscience</i> , 2.	<i>Processing Systems, NIPS ’21</i> , Red Hook, NY, USA.	395
344	Joseph B Kruskal. 1964. Multidimensional scaling by	Curran Associates Inc.	396
345	optimizing goodness of fit to a nonmetric hypothesis.	John Wu, Yonatan Belinkov, Hassan Sajjad, Nadir Dur-	397
346	<i>Psychometrika</i> , 29(1):1–27.	rani, Fahim Dalvi, and James Glass. 2020. Similar-	398
347	Brian Kulis et al. 2013. Metric learning: A sur-	ity analysis of contextual word representation mod-	399
348	vey. <i>Foundations and Trends® in Machine Learning</i> ,	els . In <i>Proceedings of the 58th Annual Meeting of</i>	400
349	5(4):287–364.	<i>the Association for Computational Linguistics</i> , pages	401
		4638–4655, Online. Association for Computational	402
		Linguistics.	403

Appendices

A The Probing Dataset

The probing dataset contains a list of sentences and their corresponding list of linguistic features. Sentences and features were generated using a custom grammar to cover central linguistic features, such as grammatical number, gender or tense, as well as confounding factors, such as word position. Table S1 shows several sentence examples, and the marking of features for each word.

To secure a clean interpretation of the relative contributions of the different features, we checked for correlations between linguistic features. Figure S5 shows the pairwise Pearson correlations among all features in the dataset.

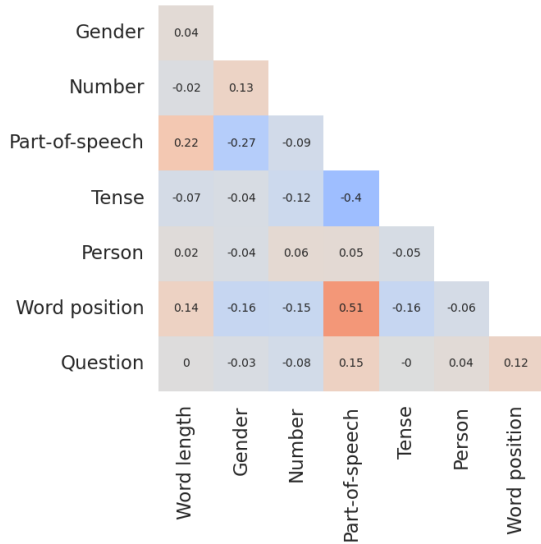


Figure S5: Pairwise Pearson correlations among all linguistic features in the probing dataset.

B Metric Learning Encoding Models (MLEMs)

Following Jalouzet et al. (2024), we provide a formal description of a Metric-Learning Encoding Model: Consider a set of N sentences, each characterized by a set of (linguistic) features \mathcal{F} . MLEMs compute two types of pairwise distances. First, *pairwise neural distances* $D^{\mathcal{N}}$ (right branch in Fig. S6), which are computed based on standard distance (e.g. Euclidean or cosine distance) between the neural responses of a set of units (e.g. a layer) for any two sentences.

Second, *pairwise feature distances* $D^{\mathcal{F},W}$ are computed as follows. First, *feature difference vectors* are computed, which indicate on which

features two sentences differ: $\Delta(s_i, s_j) = (\mathbb{1}_{f(s_i) \neq f(s_j)})_{f \in \mathcal{F}}$. Then, feature distances are computed using a standard bi-linear form parameterized by a symmetric positive definite matrix $W \in \mathbb{M}_n^+$:

$$\left(D_{ij}^{\mathcal{F},W}\right)^2 = \Delta(s_i, s_j)^T W \Delta(s_i, s_j)$$

MLEMs, as metric-learning methods, optimize W to bring the pairwise feature distances as close as possible to the neural ones, across all (i, j) pairs of stimuli:

$$W^* = \operatorname{argmin}_{W \in \mathbb{M}_n^+} \sum_{i < j} \left(\left(D_{ij}^{\mathcal{F},W}\right)^2 - \left(D_{ij}^{\mathcal{N}}\right)^2 \right)^2 + \lambda \|W\|_2^2$$

When W is assumed to be diagonal (with no interaction terms), the optimization problem simplifies to a least-squares problem, and the symmetric positive definite constraint transforms into a non-negativity constraint on the diagonal elements.

Model Training and Evaluation As in Jalouzet et al. (2024), for simplicity, we focused on the diagonal case of W and trained a standard Ridge model with a non-negativity constraint on the parameters. The regularization parameter α was optimized using nested cross-validation (CV; $\alpha \in 10^{[-4,4]}$). To facilitate α optimization across all models, target values were min-max scaled into $[0, 1]$. We evaluated the model using the Spearman correlation score ρ and report the average across CV splits. This score only assesses the similarity between the ranks of the predictions and those of the ground-truth. We chose this score as it is independent

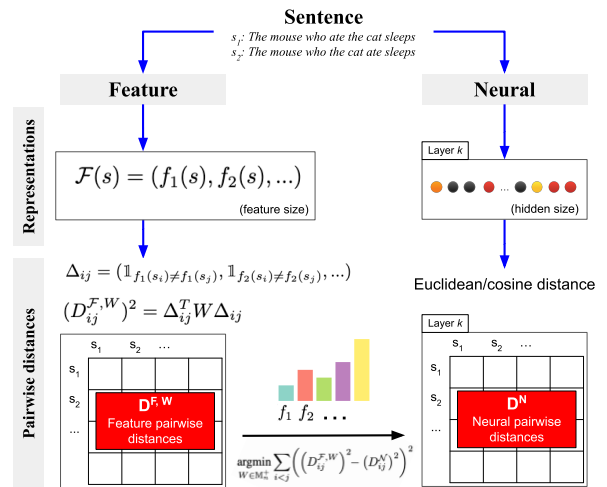


Figure S6: A Metric-Learning Encoding Model: MLEMs determine the relative importance of features by identifying the optimal alignment between distances in feature space and neural space.

word	Word length	Gender	Number	PoS	Tense	Person	Word position	Question
the	3	NaN	NaN	Det	NaN	NaN	0	False
woman	5	female	singular	Noun	NaN	3	1	False
plays	5	Nan	singular	Verb	present	3	2	False
no	2	NaN	NaN	Det	NaN	NaN	0	False
prince	6	male	singular	Noun	NaN	3	1	False
sings	5	NaN	singular	Verb	present	3	2	False
I	1	NaN	singular	Pronoun	NaN	1	0	False
vanished	8	NaN	NaN	Verb	past	NaN	1	False
do	2	NaN	singular	Auxiliary	present	NaN	0	True
you	3	NaN	NaN	Pronoun	NaN	2	1	True
sing	4	NaN	NaN	Verb	present	NaN	2	True
Mary	4	female	singular	Noun	NaN	3	0	False
fell	4	NaN	NaN	Verb	past	NaN	1	False
which	5	NaN	NaN	Wh-word	NaN	NaN	0	True
men	3	male	plural	Noun	NaN	3	1	True
sneezed	7	NaN	NaN	Verb	past	NaN	2	True

Table S1: **Examples from the probing dataset**

454 of the scale of the data (unlike the Mean Squared
455 Error) and it cannot be arbitrarily negative when
456 the estimator is very bad (unlike the coefficient of
457 determination R^2).