

Grounded Concreteness: Human-Like Concreteness Sensitivity in Vision–Language Models

Anonymous ACL submission

Abstract

Do vision–language models (VLMs) develop more human-like sensitivity to linguistic concreteness than text-only large language models (LLMs) when both are evaluated with text-only prompts? We study this question with a controlled comparison between matched Llama text backbones and their Llama Vision counterparts across multiple model scales, treating multimodal pretraining as an ablation on perceptual grounding rather than access to images at inference. We measure concreteness effects at three complementary levels: (i) output behavior, by relating question-level concreteness to QA accuracy; (ii) embedding geometry, by testing whether representations organize along a concreteness axis; and (iii) attention dynamics, by quantifying context reliance via attention-entropy measures. In addition, we elicit token-level concreteness ratings from models and evaluate alignment to human norm distributions, testing whether multimodal training yields more human-consistent judgments. Across benchmarks and scales, VLMs show larger gains on more concrete inputs, exhibit clearer concreteness-structured representations, produce ratings that better match human norms, and display systematically different attention patterns consistent with increased grounding.

1 Introduction

Human meaning is not uniformly “linguistic”: some concepts are tightly linked to perception and action (e.g., *apple*, *run*), while others are largely relational and context-dependent (e.g., *stronger*, *justice*). A long tradition in cognitive science treats *concreteness* as a graded dimension of conceptual representation, with concrete words benefiting from richer sensory codes and exhibiting robust behavioral advantages over abstract words (Paivio, 1990; Barsalou, 2008). Concreteness therefore offers a rare bridge between cognitive theory (how humans represent meaning) and computational diagnostics (how models encode and use meaning),

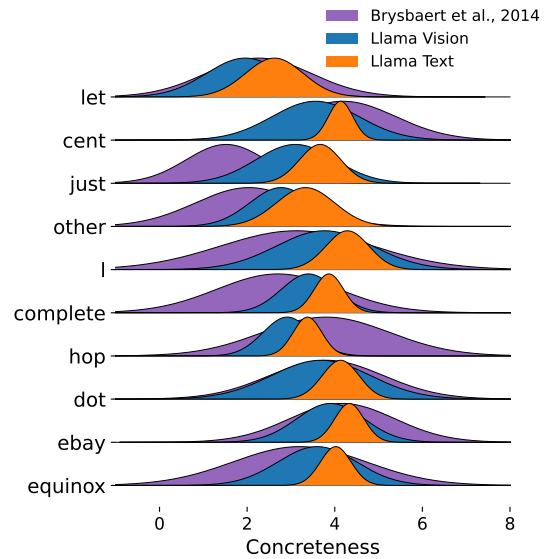


Figure 1: Comparison of concreteness rating distributions for selected words. For each word, we plot the empirical distribution of model-generated token ratings from Llama Vision (VLM) and Llama Text (LLM), alongside human norms from Brysbaert et al. (2014).

enabling measurable tests of *cognitive alignment* between humans and modern language systems (Coltheart, 1981; Brysbaert et al., 2014). More broadly, recent work in language acquisition argues that neural models can serve as hypothesis generators and testers for cognitive theories, provided we design analyses that connect internal mechanisms to behavioral signatures (Portelance, 2022a).

A central tension is that contemporary large language models (LLMs) learn from text alone, raising questions about whether “meaning” can be recovered from form without grounding (Harnad, 1990; Bender and Koller, 2020). While distributional learning can capture many semantic regularities, the absence of perceptual experience may be especially consequential for *concreteness*: in humans, concrete concepts are supported by sensorimotor simulations and imagery-like codes (Paivio, 1990;

Barsalou, 2008). Vision-language models (VLMs) offer a natural testbed for this debate. By aligning text with visual representations (e.g., CLIP-style contrastive learning and instruction-tuned VLMs), VLMs may develop more human-like, graded concreteness representations than comparably sized text-only LLMs (Radford et al., 2021; Alayrac et al., 2022; Liu et al., 2023; Touvron et al., 2023). Yet prior evidence connecting concreteness to model behavior and representations is difficult to interpret as a *vision effect*. On one hand, multi-modal/distributional semantics work links concreteness to perceptual grounding and visual consistency (Hill et al., 2014; Hessel et al., 2018; Mickus et al., 2023); on the other hand, separate lines probe concreteness as an axis in embedding spaces using text-derived features (Charbonnier and Wartena, 2019; Wartena, 2024). However, these strands rarely provide a controlled ablation that isolates the contribution of visual input, and they typically analyze a single level (task performance *or* representations) rather than jointly linking behavior, geometry, and processing dynamics.

This motivates a controlled LLM–VLM ablation in which the language backbone is held as comparable as possible and the primary difference is access to vision, allowing the contribution of visual grounding to be isolated. In this work, concreteness awareness is evaluated under this ablation by triangulating evidence across three complementary levels of analysis. First, at the **output level**, we ask whether VLM accuracy on question answering increases with the concreteness of the queried concepts. We further elicit model-produced concreteness ratings (token-level) and measure their alignment with human norms (Figure 1) (Coltheart, 1981; Brysbaert et al., 2014). Second, at the **embedding level**, we test whether token representations organize by graded concreteness by projecting token-level embeddings into a low-dimensional space (t-SNE) and measuring within-bin compactness via intra-cluster dispersion across concreteness bins. Third, at the **attention level**, we quantify contextual dependence as the entropy of each token’s self-attention distribution: if abstract meaning is more compositionally supported by surrounding context, abstract tokens should exhibit broader, higher-entropy attention, whereas concrete tokens should exhibit sharper, lower-entropy attention concentrated on fewer positions. This prediction aligns with classic context availability accounts of concreteness effects in psycholinguistics, which argue

that abstract words benefit more from supportive context for comprehension than concrete words (Schwanenflugel and Shoben, 1983; Schwanenflugel et al., 1992).

Our analyses are designed as a vision ablation for semantic grounding: holding the base language model family and scaling regime as constant as possible, we ask what changes when vision is introduced. This yields a coherent story linking classic cognitive accounts of concreteness—dual coding and grounded cognition (Paivio, 1990; Barsalou, 2008)—to measurable signatures in modern foundation models: (i) *behavioral sensitivity* to concreteness in downstream QA, (ii) *representational geometry* that recovers a concreteness ordering, and (iii) *contextual dependence* reflected by attention entropy.

Contributions. We make following contributions: (1) a controlled LLM–VLM ablation study that isolates the effect of visual grounding on concreteness awareness across matched model families and scales; (2) output-level evidence that tests concreteness sensitivity in QA and quantifies alignment between model-elicited concreteness ratings and human norms; (3) internal diagnostics connecting grounding-based accounts to model representations and processing, including concreteness-conditioned clustering in embedding space (t-SNE with intra-cluster dispersion) and attention-entropy measures of contextual dependence motivated by context-availability theories of abstract meaning.

2 Related work

2.1 Measuring concreteness in language

Concreteness is a graded psycholinguistic property that captures how directly a concept can be experienced through the senses, and it has been extensively measured through human norming studies. Classic resources such as the MRC Psycholinguistic Database provide lexical attributes including concreteness/tangibility judgments (Coltheart, 1981), and later large-scale norms substantially expand coverage and improve reliability for modern evaluation settings (Brysbaert et al., 2014).

In parallel, computational work has proposed algorithmic approximations of concreteness. Early NLP approaches connected concreteness-related cues to figurative language phenomena, using concrete vs. abstract contextual signals for literal/metaphorical sense identification (Turney et al.,

163	2011). More recent methods treat concreteness prediction as a supervised estimation problem over	and represent concreteness.	214
164	distributional features and contextual representations		
165	(Charbonnier and Wartena, 2019). A particularly	2.3 Neural models as cognitive probes of	215
166	relevant line incorporates visual information: “visual concreteness” can be operational-	language learning and processing	216
167	ized via cross-image consistency within multi-	A growing cognitive-science perspective treats	217
168	modal datasets (Hessel et al., 2018), and visually	neural networks as tools for generating and test-	218
169	grounded learning objectives can shape linguistic	ing mechanistic hypotheses about human learn-	219
170	structure and representations (Shi et al., 2019).	ing, rather than purely as engineering solutions	220
171	Given known context sensitivity (a word’s per-	(Portelance, 2022b). This includes using language	221
172	ceived concreteness can shift with discourse and	models as “psycholinguistic subjects,” evaluating	222
173	reference), human norms provide an external an-	whether model-based surprisal and state representa-	223
174	chor for evaluation, while model-produced ratings	tions predict human processing difficulty and syn-	224
175	can be treated as context-conditioned distributions	tactic expectations (Goodkind and Bicknell, 2018;	225
176	rather than fixed type-level attributes.	Futrell et al., 2019). Another line tests whether	226
177		models acquire human-relevant grammatical gen-	227
178		eralizations via targeted syntactic evaluations and	228
179	2.2 Grounding and vision language models	minimal-pair benchmarks (Linzen et al., 2016; Gu-	229
180	Grounding-based accounts of meaning empha-	lordava et al., 2018; Warstadt et al., 2020). More	230
181	size that linguistic symbols ultimately connect	recent work pushes toward developmental plausi-	231
182	to perception and action, motivating the classic	bility by constraining data and supervision (e.g.,	232
183	symbol-grounding problem (Harnad, 1990). In	BabyLM) or by studying interactive learning dy-	233
184	NLP, grounding has been studied through both	namics (Warstadt et al., 2023; Ma et al., 2025).	234
185	definitional discussions and benchmark/task	Together, these efforts motivate treating representa-	235
186	design, with critiques highlighting that “grounding”	tional properties of language models as empirical	236
187	can mean different things depending on modality,	objects for studying human cognitive constructs.	237
188	interaction, and evaluation protocol. This moti-		
189	vates evaluating grounding beyond downstream	3 Experiment Setup	238
190	success rates, using complementary diagnostics	Models The study compares matched pairs	239
191	that probe internal representations and process-	of text-only LLMs and vision-language models	240
192	ing rather than relying on task behavior alone	(VLMs) at two parameter scales. For the LLMs,	241
193	(Bisk et al., 2020a; Chandu et al., 2021; Mickus	the backbone is Meta’s Llama 3.1 family (an 8B	242
194	et al., 2023).	and a 70B Text-only Model) (Meta AI, 2024a).	243
195	Modern VLMs provide a scalable route to	For VLMs, the corresponding vision models from	244
196	grounding by learning joint representations of	the Llama Vision 3.2 family (an 11B and a 90B	245
197	text and vision. Several models are explicitly	vision LLM) (Meta AI, 2024b) were chosen that	246
198	designed to encourage fine-grained grounding	utilize the Llama 3.1 text-only models as a	247
199	via cross-modal alignment objectives (e.g.,	backbone. The text-only models are fitted with	248
200	word/phrase-region alignment) (Tan and Bansal,	a vision adapter and then trained on a multi-	249
201	2019; Chen et al., 2020), and to evaluate	modal dataset to create the vision models. We	250
202	grounded lexical acquisition beyond standard	refer Appendix A for details. Unless otherwise	251
203	downstream transfer (Ma et al., 2023). Large-	stated, evaluation uses text-only prompts	252
204	scale contrastive and generative pretraining	(no images), so the LLM–VLM comparison	253
205	has produced general-purpose models that	functions as an ablation on <i>access to visual</i>	254
206	align linguistic descriptions with visual	<i>supervision during training</i> rather than access	255
207	features, supporting transfer to many multi-	to images at inference.	
208	modal tasks (Radford et al., 2021; Alayrac et	Measuring concreteness Token-level concre-	256
209	al., 2022; Liu et al., 2023).	teness $C(w)$ is obtained from the human	257
210	These strands motivate treating vision as a	ratings in the 40k English words (Brysbaert	258
211	causal factor that can strengthen concreteness	et al., 2014) (40K). For each word that	259
212	awareness. If concrete concepts are more	appears in 40K, its concreteness score is	260
213	consistently tied to perceptual regularities	set to the corresponding 40K mean rating.	261
	(i.e., they have more stable visual corre-	For out-of-vocabulary proper nouns	262
	lates), then adding visual supervision	(e.g., named entities) that are not covered	
	should preferentially benefit how models	by 40K,	
	recognize		

the score is set to the maximum concreteness value on the 40K scale (5). For function words without a clear concreteness interpretation (e.g., articles and prepositions) that are also absent from 40K, the score is set to 0. Sentence-level concreteness for an input string x with word tokens $w_{1:n}$ is the mean of token scores:

$$C(x) = \frac{1}{n} \sum_{i=1}^n c(w_i). \quad (1)$$

For subword-tokenized model inputs, word-level scores are propagated to constituent sub-tokens to enable tokenwise analyses.

Text datasets Evaluation uses standard text-only QA benchmarks covering diverse reasoning demands and a broader range of question concreteness: ARC-Easy and ARC-Challenge for grade-school science multiple-choice questions (Clark et al., 2018); BoolQ for naturally occurring yes/no questions (Clark et al., 2019); Winogrande for adversarial pronoun/coreference resolution (Sakaguchi et al., 2020); CommonsenseQA for commonsense multiple-choice QA (Talmor et al., 2019); Social IQA for reasoning about social interactions and implications (Sap et al., 2019) and PIQA for physical commonsense reasoning (Bisk et al., 2020b). Performance is measured by accuracy under a unified prompting format. We refer Appendix B and C for details on datasets and prompts.

3.1 Research questions and methods

This section describes how each hypothesis is operationalized and tested. All analyses are conducted for both model scales to assess scaling effects.

Does the VLM outperform the LLM on QA questions as question concreteness increases?

For each benchmark dataset, each question is scored as correct or incorrect under a unified prompting format. To summarize performance as a function of concreteness, sentence-level concreteness scores are pooled across all datasets and discretized into six equal-width bins of size 0.6, spanning [1.8, 4.8] (the observed range is [1.96, 4.67]). For each bin, accuracy is computed as the mean correctness over questions whose sentence concreteness falls in that interval. In addition, the bin-wise accuracy gap between the VLM and its matched LLM is reported, $\Delta\text{Acc} = \text{Acc}_{\text{VLM}} - \text{Acc}_{\text{LLM}}$, to quantify where vision provides an advantage. We

hypothesize that ΔAcc is expected to be larger in higher-concreteness bins, indicating that the VLM is relatively more robust on concrete questions than the text-only model.

Do VLM token representations exhibit tighter within-concreteness clusters than LLM token representations?

Each word (w) covered by 40K is rounded to a discrete concreteness bin $b(w) \in \{1, \dots, 5\}$. For each model, we extract last-layer contextual representations and average over occurrences to obtain a type vector $\bar{\mathbf{h}}(w)$. To measure within-bin dispersion, we fit $\bar{\mathbf{h}}(w)$ with 2D t-SNE, yielding $\mathbf{z}(w)$. Within this t-SNE space, dispersion is measured as the mean pairwise cosine distance among tokens with the same label, where lower values indicate more compact clusters:

$$D = \mathbb{E}_{\ell} \mathbb{E}_{w \neq w' \sim \mathcal{W}_{\ell}} [1 - \cos(\mathbf{z}(w), \mathbf{z}(w'))] \quad (2)$$

where \mathcal{W}_{ℓ} represents the discretized concreteness bin. Lower D indicates tighter within-concreteness clusters.

Do abstract tokens exhibit higher-entropy attention distributions than concrete tokens, and is this abstract-concrete separation sharper in VLMs?

For each layer ℓ and head h , self-attention weights follow the standard Transformer definition (Vaswani et al., 2017):

$$\mathbf{A}^{(\ell,h)} = \text{softmax}\left(\frac{\mathbf{Q}^{(\ell,h)}\mathbf{K}^{(\ell,h)\top}}{\sqrt{d_k}}\right), \quad (3)$$

where $\mathbf{A}_{i,j}^{(\ell,h)}$ is the attention paid by token i to token j and forms a probability distribution over j due to softmax normalization. For each token i , attention entropy is computed as:

$$H^{(\ell,h)}(i) = - \sum_j \mathbf{A}_{i,j}^{(\ell,h)} \log \mathbf{A}_{i,j}^{(\ell,h)}. \quad (3)$$

Entropy is then averaged across heads for each layer to obtain a per-token entropy score. At each layer, we test the association between token concreteness $c(w_i)$ and attention entropy via Pearson’s r . We expect a *negative* correlation ($r < 0$): abstract tokens should exhibit *higher* attention entropy (more diffuse context integration), whereas concrete tokens should exhibit *lower* entropy (more focused attention), consistent with concreteness effects in comprehension (Schwanenflugel and Shoben, 1983; Schwanenflugel et al., 1992). Moreover, we predict this effect is stronger in VLMs than LLMs (i.e., more negative r in VLMs), reflecting a sharper abstract-concrete separation.

Do model generated concreteness judgments align better with human norms for VLMs?

Each model is prompted to output a concreteness rating for every word token in each question on the 40K scale. We refer Appendix C for prompt details. To elicit reliable concreteness judgments and enforce a consistent output format, we use the instruct variants of the larger models (70B text-only and 90B vision-language). Because the same word type can appear in multiple contexts, each word w induces an empirical distribution over ratings under a model m :

$$p_m(r | w) \propto \sum_{x \in \mathcal{X}(w)} \mathbf{1}[\hat{r}_m(w, x) = r] \quad (4)$$

where $\mathcal{X}(w)$ are contexts containing w and $\hat{r}_m(w, x)$ is the model-produced rating for token w in question x . We construct an analogous human distribution $p_H(r | w)$ from the 40K annotations and quantify human-model agreement with the symmetric KL divergence:

$$D_{\text{KL}}(w) = \frac{1}{2} [\text{KL}(p_m \| p_H) + \text{KL}(p_H \| p_m)] \quad (5)$$

where smaller $D_{\text{KL}}(w)$ indicates better alignment. We expect (i) $D_{\text{KL}}(w)$ decreases as human concreteness increases, and (ii) this decrease is steeper for VLMs than for LLMs. To test the trend, we bin words by human concreteness in 0.5-wide bins and regress binned D_{SKL} on bin center, reporting slope, R^2 , and p -value.

4 Analysis and discussion

VLM outperform their LLM counterpart on QA questions and the gaps widen on more concrete questions. Figure 2 shows that, across all concreteness bins, VLMs consistently outperform their text-only counterparts at both scales. Across all datasets, for the smaller pair, accuracy increases from 68.0% (LLM) to 77.5% (VLM), and for the larger pair from 78.8% (LLM) to 85.5% (VLM). Since each VLM is trained from the same model family as its LLM counterpart, these gains indicate that multimodal training transfers to improved textual QA. We refer Appendix D for additional per-dataset results.

Crucially, the advantage is not uniform across question types: the bottom row of Figure 2 shows that the VLM-LLM gap increases with question concreteness in both scales, with the strongest separation in the most concrete bin. This pattern suggests that visual grounding disproportionately benefits questions whose successful resolution depends

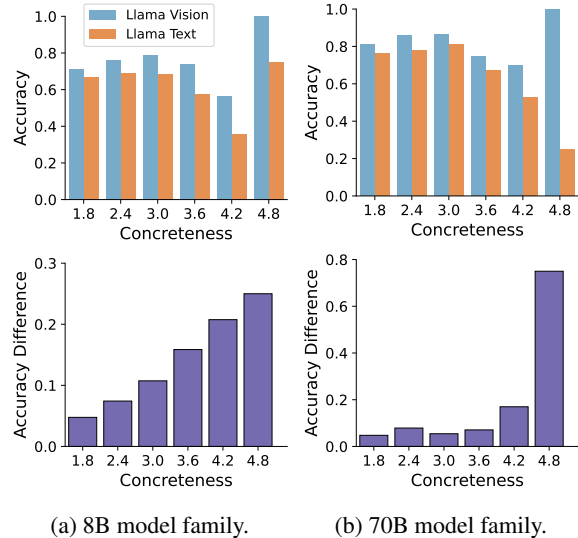


Figure 2: Top row: accuracy by question concreteness for Llama Text vs. Llama Vision. Bottom row: the VLM-LLM accuracy gap.

on perceptible entities, attributes, and events (e.g., shape, material, spatial relations), consistent with grounded accounts in which perceptual experience provides an additional scaffold for semantic representations (Harnad, 1990; Barsalou, 2008; Bisk et al., 2020a). A plausible mechanism is that vision-language training strengthens the association between concrete lexical items and perceptually anchored image features (e.g., object properties and spatial configurations), making the relevant evidence easier to retrieve and compose when answering concrete questions. In other words, the VLM’s gains appear concentrated where the QA signal can be supported by grounded semantics rather than purely symbolic co-occurrence.

At the same time, the effect is smaller (and sometimes flatter) in lower-concreteness bins, where performance may depend more on abstract relations, discourse-level inference, or world knowledge not directly supported by perceptual grounding. Moreover, abstract language tends to be more polysemous and context-dependent, which can reduce the benefit of any single additional modality. Together, these results suggest that multimodal training provides an asymmetric benefit: it reliably improves QA overall, but disproportionately improves the processing and use of concrete concepts, which we further probe via representation geometry and attention diagnostics in subsequent experiments.

VLM token representations form tighter within-concreteness clusters than LLMs. Given that

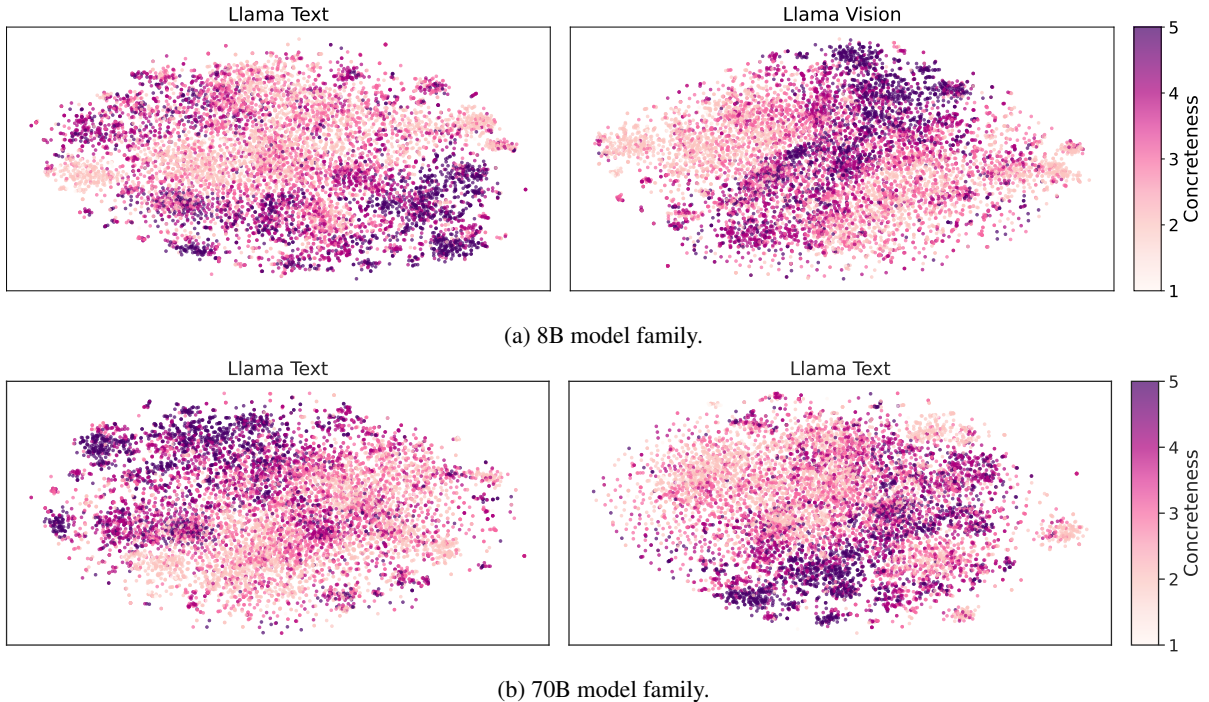


Figure 3: t-SNE of average last-layer token representations for Llama Text vs. Llama Vision, colored by human concreteness.

VLMs outperform their text-only counterparts on QA, we ask whether vision-text training also reshapes the *geometry* of token representations along a graded concreteness dimension. Figure 3a provides a qualitative view: in both the smaller and larger models, the VLM embedding shows a visibly more contiguous band of highly concrete tokens (darker points), whereas the text-only LLM distributes high-concreteness tokens more diffusely across the map. This pattern suggests that visual grounding encourages representations of perceptually grounded words to occupy a more coherent subregion of the space, consistent with grounded accounts of meaning and the symbol-grounding perspective (Harnad, 1990; Barsalou, 2008; Bisk et al., 2020a).

We quantify this effect using within-bin intra-cluster dispersion (Eq. 2) computed in the t-SNE space. Table 1 shows that VLMs achieve lower dispersion than LLMs at *every* concreteness level in both families. The effect is largest for the most concrete bin ($c=5.0$): dispersion drops from $0.76 \rightarrow 0.66$ in the 8B family and from $0.87 \rightarrow 0.77$ in the 70B family. Notably, dispersion is highest in the mid-concreteness range (roughly $c \in [2, 4]$) and drops sharply for the most concrete words, which is compatible with the idea that mid-range words are more heterogeneous (e.g., broader senses or

Conc.	8B model family		70B model family	
	Text-only ↓	Vision ↓	Text-only ↓	Vision ↓
1.0	0.87	0.75	0.93	0.82
2.0	0.94	0.87	0.98	0.94
3.0	0.98	0.96	1.00	0.99
4.0	0.99	0.96	1.00	0.99
5.0	0.76	0.66	0.87	0.77

Table 1: Within-concreteness cluster dispersion (mean pairwise cosine distance in 2D t-SNE; lower is tighter).

mixed perceptual/abstract usage) while highly concrete words admit more stable, visually grounded semantics.

Taken together, the qualitative structure in Figure 3a and the consistent quantitative reductions in Table 1 support our hypothesis that VLM representations encode graded concreteness more cleanly than LLMs. This geometry offers a representational account of the concreteness-dependent QA gains. If highly concrete word types occupy a tighter region of the space, their representations are more consistent across contexts, reducing the need for context-dependent disambiguation and making grounded attributes easier to retrieve and compose. In turn, this should improve robustness on questions that hinge on perceptual properties (e.g., materials, shapes, spatial relations).

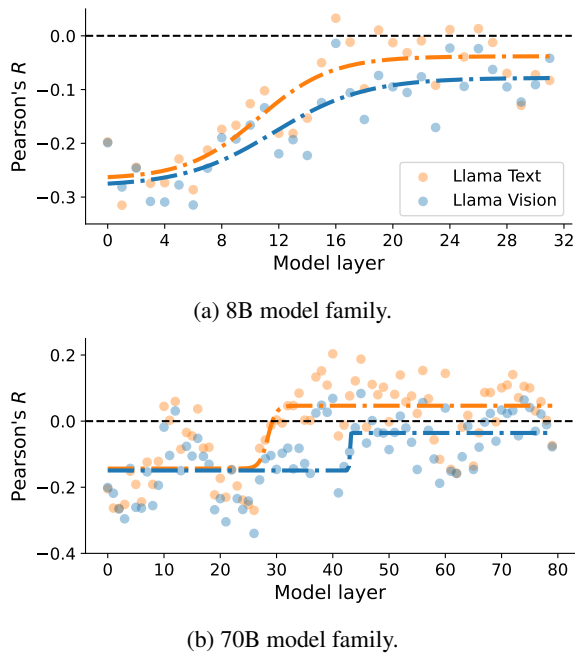


Figure 4: Layerwise Pearson’s r between token concreteness and head-averaged attention entropy. Colored dash lines are sigmoid fitted curves.

Concrete tokens show lower attention entropy, with a stronger effect in VLMs. Motivated by the representational results showing tighter within-concreteness clustering in VLMs, we next test whether models also differ in how much context they integrate for abstract versus concrete words. Prior psycholinguistic work suggests that concrete words tend to be easier to interpret and rely less on contextual support than abstract words, which are more context-dependent (Schwanenflugel and Shoben, 1983; Schwanenflugel et al., 1992). We operationalize contextual reliance using attention entropy (Eq. 3): higher entropy corresponds to more diffuse attention over many tokens, while lower entropy reflects more concentrated attention.

Figure 4 plots, for each layer, Pearson’s r between token concreteness and token attention entropy (head-averaged), for both model scales. We report full Pearson’s R values and their p -values in Appendix E. Across most layers, correlations are negative (or near zero), indicating that more concrete tokens exhibit lower attention entropy (i.e., they attend more selectively) whereas abstract tokens show higher entropy, consistent with the hypothesis that abstract meaning requires broader contextual integration. The effect is strongest in earlier-to-mid layers: both model families display more negative correlations in roughly the first half

of the network, followed by a gradual attenuation toward later layers. This layerwise result suggests that context concreteness sensitivity is primarily expressed early in processing, while later layers may shift toward task-level integration that is less directly tied to lexical concreteness.

Importantly, VLMs show a consistently stronger negative correlation than their text-only counterparts. Averaging Pearson’s r across layers yields $r = -0.12$ (LLM) vs. $r = -0.16$ (VLM) for the smaller models, and $r = -0.02$ (LLM) vs. $r = -0.10$ (VLM) for the larger models, indicating a sharper abstract–concrete separation in VLM attention behavior. One interpretation is that vision-text training provides an additional grounding signal that stabilizes the representations of concrete words, allowing the model to resolve them with more focused attention (lower entropy) and reducing the need to distribute attention broadly across other tokens.

VLM concreteness ratings align more closely with humans than LLM as concreteness increases. Finally, if the performance, representation geometry, and attention analyses reflect a human-like graded concreteness sensitivity, we should also observe human-aligned concreteness judgments from the models. Figure 5 plots human concreteness against human–model agreement measured by symmetric KL divergence (lower is better). Overall, the VLM exhibits lower divergence than the LLM (mean D_{KL} : 9.4 vs. 10.1), indicating closer alignment to human norms from (Brysbaert et al., 2014).

More importantly, agreement improves with concreteness: as human ratings increase, D_{KL} decreases for both models, but the trend is substantially sharper for the VLM. A linear fit on binned concreteness shows that the VLM’s alignment increases reliably with concreteness (slope = -1.810 , $R^2 = 0.857$, $p < 0.001$), whereas the LLM trend is weaker and not statistically reliable (slope = -1.048 , $R^2 = 0.328$, $p > 0.1$). This suggests that vision-text training does not merely shift ratings globally, but preferentially calibrates judgments for perceptually grounded words—precisely where vision provides an additional supervisory signal.

Because concreteness is an interpretable, human-normed semantic axis, improved human–model alignment makes model behavior easier to diagnose: it supports using concreteness as a principled factor for error analysis (e.g., when failures concen-

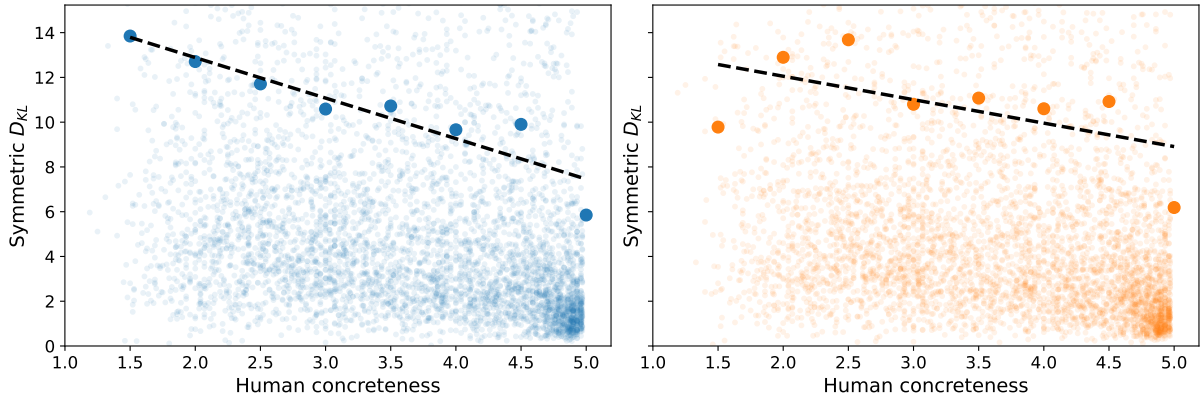


Figure 5: Human-model alignment of token-level concreteness judgments. Each point is a word with human concreteness on the x-axis and symmetric KL divergence between the model’s rating distribution and the human (40K) distribution on the y-axis (lower is better). Larger dots show bin averages; dashed lines are linear fits over bins. The VLM (left) exhibits lower divergence and a steeper decrease in divergence with concreteness than the matched text-only LLM (right).

trate in abstract language) and as an interpretable control variable when comparing model families, aligning with calls for more rigorous, task-relevant interpretability evaluations (Doshi-Velez and Kim, 2017; Guidotti et al., 2019). Mechanistically, concreteness alignment provides a concrete target for circuit-level analysis: one can localize where concreteness enters computation (layers/heads/MLP features) and test causal interventions, complementing transformer reverse-engineering frameworks (Elhage et al., 2021; Geva et al., 2022).

Model concreteness behaviour carries through scale. Across all analyses, the concreteness effects observed in the 8B model family qualitatively mirror in the 70B family: VLMs show larger gains on concrete QA, tighter within-concreteness clustering, more negative concreteness-entropy correlations, and stronger human-aligned rating trends. This consistency suggests that concreteness organization is not specific to small models, but a stable property that persists under scaling, with vision-text training providing an additive grounding signal rather than a scale-specific artifact.

5 Conclusion

We presented a controlled test of whether visual supervision during training induces more human-like concreteness sensitivity in foundation models. Holding the language backbone family and scaling regime as constant as possible and using text-only evaluation prompts, we compared matched Llama 3.1 LLMs against their Llama Vision 3.2 counterparts, treating the LLM-VLM contrast as an

ablation on access to perceptual grounding rather than access to images at inference. Across diverse QA benchmarks, VLMs achieved higher accuracy overall and, crucially, showed larger improvements on questions with higher concreteness, consistent with grounded cognition and dual-coding accounts in which perceptual experience disproportionately supports concrete semantics (Paivio, 1990; Barsalou, 2008; Harnad, 1990). Internal analyses converged on a coherent mechanistic picture: VLM token representations formed tighter within-concreteness clusters in low-dimensional projections, suggesting more stable type-level semantics for highly concrete words; and attention-entropy diagnostics indicated a sharper abstract-concrete separation in contextual reliance, aligning with psycholinguistic theories that abstract meaning draws more heavily on supportive context (Schwanenflugel et al., 1992; Schwanenflugel and Shoben, 1983). Finally, elicited token-level concreteness ratings agreed more closely with human norms in VLMs, with a stronger improvement in human-model alignment as concreteness increased, indicating that vision-text training preferentially calibrates judgments precisely where perceptual supervision is informative.

Beyond the performance gap, our results position concreteness as a principled, interpretable axis for comparing model families and diagnosing grounding-related behavior. This provides a useful bridge between cognitive constructs and modern interpretability practice: concreteness offers a measurable target for localizing where grounded semantic information enters computation.

625 Limitations

626 **Single model family.** Our controlled comparison
627 centers on one matched LLM/VLM family, im-
628 proving internal validity but limiting generality. A
629 direct next replication is to repeat the same pipeline
630 on additional matched pairs, such as multilingual
631 families Qwen/Qwen-VL to test whether concrete-
632 ness effects and human-alignment patterns persist
633 across languages and training recipes (Bai et al.,
634 2023; Wang et al., 2024).

635 **Token frequency as a confound.** Some apparent
636 concreteness effects may be partly explained by
637 lexical frequency/contextual diversity rather than
638 grounding, but we cannot access true pretraining
639 token counts. To reduce this confound, we can
640 add frequency proxies (tokenizer-matched counts
641 from large open corpora, plus average surprisal on
642 held-out text) as covariates, or frequency-match
643 concrete vs. abstract item sets before running the
644 main regressions.

645 **Developmental trajectory.** We only analyze fi-
646 nal checkpoints, so we cannot determine *when* con-
647 creteness sensitivity and human-alignment emerge
648 or how this depends on scale. An immediate follow-
649 up is a developmental-style study over intermedi-
650 ate checkpoints (or staged training) to track the
651 concreteness–accuracy slope and internal separa-
652 bility over training time for each model scale.

653 References

654 Jean-Baptiste Alayrac, Jeffrey Donahue, Pauline Luc,
655 Antoine Miech, Iain Barr, Yana Hasson, and 1 others.
656 2022. [Flamingo: a visual language model for few-
657 shot learning](#). In *Advances in Neural Information
658 Processing Systems*.

659 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,
660 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,
661 and Jingren Zhou. 2023. [Qwen-vl: A versatile vision-
662 language model for understanding, localization, text
663 reading, and beyond](#). *Preprint*, arXiv:2308.12966.

664 Lawrence W. Barsalou. 2008. [Grounded cognition](#). *An-
665 nual Review of Psychology*, 59:617–645.

666 Emily M. Bender and Alexander Koller. 2020. [Climbing
667 towards NLU: On meaning, form, and understanding
668 in the age of data](#). In *Proceedings of the 58th Annual
669 Meeting of the Association for Computational Lin-
670 guistics*. Association for Computational Linguistics.

671 Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob
672 Andreas, Yoshua Bengio, Joyce Chai, Mirella Lap-
673 ata, Angeliki Lazaridou, Jonathan May, Aleksandr

Nisnevich, Nicolas Pinto, and Joseph Turian. 2020a. [Experience grounds language](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics. 674
675
676
677
678

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jian-
feng Gao, and Yejin Choi. 2020b. [PIQA: Reasoning
about physical commonsense in natural language](#). In
*Proceedings of the AAAI Conference on Artificial
Intelligence*. 679
680
681
682
683

Marc Brysbaert, Amy Beth Warriner, and Victor Ku-
perman. 2014. [Concreteness ratings for 40 thousand
generally known English word lemmas](#). *Behavior
Research Methods*, 46(3):904–911. 684
685
686
687

Khyathi Chandu, Siva Reddy, Alan W. Black, Yu-
lia Tsvetkov, and Eric Nyberg. 2021. [Grounding
“grounding” in NLP](#). In *Findings of the Association
for Computational Linguistics: ACL-IJCNLP 2021*.
Association for Computational Linguistics. 688
689
690
691
692

Jean Charbonnier and Christian Wartena. 2019. [Pre-
dicting word concreteness and imagery](#). In *Proceed-
ings of the 13th International Conference on Com-
putational Semantics - Long Papers*, pages 176–187,
Gothenburg, Sweden. Association for Computational
Linguistics. 693
694
695
696
697
698

Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed
El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and
Jingjing Liu. 2020. [Uniter: Universal image-text rep-
resentation learning](#). In *Proceedings of the European
Conference on Computer Vision (ECCV)*. 699
700
701
702
703

Christopher Clark, Kenton Lee, Ming-Wei Chang,
Tom Kwiatkowski, Michael Collins, and Kristina
Toutanova. 2019. [BoolQ: Exploring the surprising
difficulty of natural yes/no questions](#). In *Proceedings
of the 2019 Conference of the North American Chap-
ter of the Association for Computational Linguistics:
Human Language Technologies, Volume 1 (Long and
Short Papers)*, pages 2924–2936, Minneapolis, Min-
nesota. Association for Computational Linguistics. 704
705
706
707
708
709
710
711
712

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,
Ashish Sabharwal, Carissa Schoenick, and Oyvind
Tafjord. 2018. [Think you have solved question an-
swering? try ARC, the AI2 reasoning challenge](#).
arXiv preprint arXiv:1803.05457. 713
714
715
716
717

Max Coltheart. 1981. The MRC psycholinguistic
database. *Quarterly Journal of Experimental Psy-
chology*. 718
719
720

Finale Doshi-Velez and Been Kim. 2017. [Towards a
rigorous science of interpretable machine learning](#).
arXiv preprint arXiv:1702.08608. 721
722
723

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom
Henighan, Nicholas Joseph, Ben Mann, Amanda
Askell, Yuntao Bai, Anna Chen, Tom Conerly,
Nova DasSarma, Dawn Drain, Deep Ganguli,
Zac Hatfield-Dodds, Danny Hernandez, Andy
Jones, Jackson Kernion, Liane Lovitt, Kamal 724
725
726
727
728
729

843		900
844		901
845		902
846		903
847		904
848		905
849		906
850		907
851		908
852		909
853		910
854		911
855		912
856		913
857		914
858		915
859		916
860		917
861		918
862		919
863		920
864		921
865		922
866		923
867		924
868		925
869		926
870		927
871		928
872		929
873		930
874		931
875		932
876		933
877		934
878		935
879		936
880		937
881		938
882		939
883		940
884		941
885		942
886		943
887		944
888		945
889		946
890		947
891		948
892		949
893		950
894		
895		
896		
897		
898		
899		

Model Family	Size	Layers	Hidden Dim	Attn Heads	KV Heads
Llama 3.1 (Text)	8B	32	4096	32	8
Llama 3.2 (Vision)	11B	32	4096	32	8
Llama 3.1 (Text)	70B	80	8192	64	8
Llama 3.2 (Vision)	90B	80	8192	64	8

Table 2: Architectural specifications for the language backbones used in this study. The VLM variants inherit these text specifications and add vision encoder parameters.

WinoGrande, CommonsenseQA, PIQA), we report results on the canonical validation/development split. For ARC, we use the test split. All datasets and model checkpoints are obtained from the public Hugging Face Hub (via the datasets and transformers libraries). The benchmarks consist of generic multiple-choice questions and do not require any user-provided inputs. We do not collect, store, or process personally identifying information. Accordingly, our experiments pose minimal risk of identity disclosure, and we report only aggregate accuracy metrics without releasing any per-example outputs that could contain sensitive content.

C Prompts

QA Evaluation (Non-Instruct Models) For the base (non-instruct) models, we utilized zero-shot prompting templates tailored to the task format. For multiple-choice datasets (ARC, CommonsenseQA, PIQA, WinoGrande), we used a standard completion format that lists options and prompts for an immediate answer:

```
You are a helpful assistant.
Answer the question immediately
with just the option letter (A,
B, C, D) or number.
Question: {question}
Options:
A. {choice_a}
B. {choice_b}
C. {choice_c}
D. {choice_d}
Answer:
```

For the reading comprehension dataset (BoolQ), we employed a template that conditions the binary response on the provided passage:

```
Read the following passage and
answer the question with Yes or
No.
```

```
Passage: {passage}
```

```
Questions: {question}
```

```
Answer:
```

Concreteness Ratings To elicit concreteness ratings from the the large models (Llama 3.1 70B and Llama 3.2 90B Vision), we used the following prompt to ensure the output is aligned with the 1-7 MRC scale.

```
You are a psycholinguistics
expert. Your task is to rate the
'concreteness' of every content
word in the following text on a
scale from 1.0 (very abstract)
to 7.0 (very concrete/tangible).
```

```
Defintion:
```

```
- Concrete words refer to
things you can perceive directly
with your senses (touch, see,
hear, smell). Examples: 'apple',
'chair', 'scream'.
```

```
- Abstract words refer to
concepts, ideas, or emotions that
cannot be directly perceived.
Examples: 'freedom', 'justice',
'infinity'.
```

```
Input Text: "{text}"
```

```
Return your analysis strictly
as a JSON list of objects,
where each object has 'word' and
'score'. Ignore stop words (the,
a, is, etc.).
```

```
Example format:
```

```
[
{"word": "apple", "score": 6.2},
{"word": "freedom", "score":
2.77}
]
```

```
JSON Output:
```

Dataset	Domain	# Questions	Avg. Length
ARC-Easy	Grade-school Science	2,376	~39 words
ARC-Challenge	Grade-school Science (Hard)	1,172	~47 words
BoolQ	Reading Comprehension (Yes/No)	3,270	~120 words [†]
WinoGrande	Commonsense (Coreference)	1,267	~32 words
CommonsenseQA	General Commonsense	1,221	~26 words
PIQA	Physical Interaction	1,000	~48 words
SIQA	Social Commonsense	1,954	~36 words

Table 3: Summary of evaluation datasets. [†]Includes passage length.

D Detailed Results

Table 4 presents the raw accuracy scores for each model across all five datasets.

E Attention Entropy Analysis

In the following tables, we present the raw correlation values from the Attention Entropy analysis for both model categories. Tables 5–7 contain the Pearson correlation (R) and Significance (p) with significance levels: $**p < 0.05$, $***p < 0.01$.

Dataset	Small Category		Large Category	
	Llama 3.1 8B	Llama 3.2 11B	Llama 3.1 70B	Llama 3.2 90B
ARC-Easy	86.45%	89.23%	93.10%	93.90%
ARC-Challenge	65.70%	78.41%	89.08%	91.64%
BoolQ	76.18%	82.02%	85.02%	90.31%
WinoGrande	52.41%	63.30%	69.93%	78.85%
CommonsenseQA	61.51%	71.58%	76.17%	77.15%
PIQA	35.30%	71.90%	54.00%	74.50%
SIQA	64.43%	71.19%	65.10%	78.71%
Average				

Table 4: Per-dataset accuracy (%) for all evaluated models. The VLM variants generally perform comparable to or better than their text-only counterparts, despite receiving no visual input during this evaluation.

Table 5: Large Scale Models (Part 1, Layers 0–31)

Layer	Llama 3.1 70B		Llama 3.2 90B	
	<i>R</i>	<i>p</i>	<i>R</i>	<i>p</i>
0	-0.20	***	-0.20	***
1	-0.26	***	-0.22	***
2	-0.26	***	-0.27	***
3	-0.25	***	-0.30	***
4	-0.15	***	-0.14	***
5	-0.19	***	-0.26	***
6	-0.24	***	-0.26	***
7	-0.12	***	-0.15	***
8	-0.19	***	-0.26	***
9	-0.12	***	-0.19	***
10	0.04	***	-0.02	0.09
11	0.00	0.82	-0.10	***
12	0.06	***	0.03	***
13	-0.08	***	-0.15	***
14	-0.03	***	-0.08	***
15	-0.04	***	-0.11	***
16	0.04	***	-0.05	***
17	-0.08	***	-0.11	***
18	-0.08	***	-0.13	***
19	-0.22	***	-0.27	***
20	-0.17	***	-0.23	***
21	-0.23	***	-0.30	***
22	-0.15	***	-0.15	***
23	-0.20	***	-0.23	***
24	-0.24	***	-0.27	***
25	-0.25	***	-0.24	***
26	-0.27	***	-0.34	***
27	-0.08	***	-0.18	***
28	-0.06	***	-0.11	***
29	-0.01	0.21	-0.10	***
30	0.00	0.74	-0.15	***
31	-0.01	0.58	-0.10	***

Table 6: Large Scale Models (Part 2, Layers 32–79)

Layer	Llama 3.1 70B		Llama 3.2 90B	
	<i>R</i>	<i>p</i>	<i>R</i>	<i>p</i>
32	0.05	***	-0.08	***
33	0.05	***	-0.14	***
34	0.08	***	-0.08	***
35	0.00	0.78	-0.13	***
36	0.00	0.84	-0.16	***
37	0.13	***	0.02	**
38	0.15	***	0.05	***
39	0.11	***	0.03	**
40	0.20	***	0.07	***
41	-0.05	***	-0.22	***
42	-0.01	0.27	-0.14	***
43	0.08	***	-0.09	***
44	0.12	***	-0.02	0.06
45	0.19	***	0.08	***
46	-0.02	0.09	-0.07	***
47	0.11	***	0.00	0.88
48	0.08	***	-0.03	***
49	0.10	***	-0.04	***
50	-0.00	0.87	-0.09	***
51	0.08	***	-0.03	***
52	0.14	***	0.02	0.15
53	0.06	***	-0.02	0.06
54	0.08	***	-0.06	***
55	0.02	*	-0.15	***
56	0.15	***	0.06	***
57	0.07	***	-0.03	***
58	-0.06	***	-0.11	***
59	-0.12	***	-0.19	***
60	0.14	***	0.04	***
61	-0.15	***	-0.15	***
62	-0.16	***	-0.16	***
63	-0.02	0.09	-0.11	***
64	0.02	0.11	-0.02	0.16
65	-0.14	***	-0.15	***
66	-0.04	***	-0.08	***
67	0.09	***	-0.00	0.78
68	0.09	***	0.02	*
69	-0.00	0.83	-0.06	***
70	0.10	***	0.03	***
71	0.11	***	0.02	*
72	0.14	***	0.03	***
73	0.07	***	-0.04	***
74	0.11	***	0.06	***
75	0.10	***	0.04	***
76	0.03	***	-0.03	***
77	0.06	***	0.02	*
78	0.00	0.82	-0.01	0.31
79	-0.07	***	-0.08	***

Table 7: Small Scale Models (Layers 0–31)

Layer	Llama 3.1 8B		Llama 3.2 11B	
	<i>R</i>	<i>p</i>	<i>R</i>	<i>p</i>
0	-0.20	***	-0.20	***
1	-0.32	***	-0.28	***
2	-0.24	***	-0.25	***
3	-0.27	***	-0.31	***
4	-0.27	***	-0.31	***
5	-0.23	***	-0.28	***
6	-0.29	***	-0.31	***
7	-0.21	***	-0.25	***
8	-0.17	***	-0.19	***
9	-0.17	***	-0.19	***
10	-0.13	***	-0.17	***
11	-0.10	***	-0.13	***
12	-0.18	***	-0.22	***
13	-0.18	***	-0.19	***
14	-0.15	***	-0.22	***
15	-0.05	***	-0.12	***
16	0.03	***	-0.01	0.19
17	-0.01	0.27	-0.11	***
18	-0.10	***	-0.16	***
19	0.01	0.32	-0.07	***
20	-0.01	0.24	-0.09	***
21	-0.03	***	-0.11	***
22	-0.01	0.38	-0.08	***
23	-0.09	***	-0.17	***
24	0.01	0.28	-0.02	*
25	-0.04	***	-0.09	***
26	0.01	0.22	-0.02	**
27	-0.01	0.25	-0.06	***
28	-0.07	***	-0.10	***
29	-0.13	***	-0.12	***
30	-0.07	***	-0.09	***
31	-0.08	***	-0.04	***