# Do Multilingual Transformers Encode Paninian Grammatical Relations? A Layer-wise Probing Study

**Akshit Kumar, Dipti Sharma, Parameswari Krishnamurthy**

LTRC, International Institute of Information Technology, Hyderabad, India

`akshit.kumar@research.iiit.ac.in`

`{dipti, param.krishna}@iiit.ac.in`

## Abstract

Large multilingual transformers such as XLM-RoBERTa achieve impressive performance on diverse NLP benchmarks, but understanding how they internally encode grammatical information remains challenging. This study investigates the encoding of syntactic and morphological information derived from the Paninian grammatical framework—specifically designed for morphologically rich Indian languages—across model layers. Using diagnostic probing, we analyze the hidden representations of frozen XLM-RoBERTa-base, mBERT, and IndicBERT-v2 models across seven Indian languages (Hindi, Kannada, Malayalam, Marathi, Telugu, Urdu, Bengali). Probes are trained to predict Paninian dependency relations (by edge probing) and essential morphosyntactic features (UPOS tags, Vibhakti markers). We find that syntactic structure (dependencies) is primarily encoded in the middle-to-upper-middle layers (layers 6–9), while lexical features peak slightly earlier. Although the general layer-wise trends are shared across models, significant variations in absolute probing performance reflect differences in model capacity, pre-training data, and language-specific characteristics. These findings shed light on how theory-specific grammatical information emerges implicitly within multilingual transformer representations trained largely on unstructured raw text.

## 1 Introduction

Multilingual pre-trained transformer models such as XLM-RoBERTa (Conneau et al., 2019) have become foundational in natural language processing (NLP), achieving remarkable cross-lingual transfer capabilities. However, their internal representations remain only partially understood. Identifying where and how linguistic structures are encoded across the layers of these transformers improves their interpretability, efficiency, and adaptability across diverse languages and linguistic tasks (Belinkov and Glass, 2019).

A fundamental question is whether multilingual transformer models implicitly encode theory-specific grammatical information despite not being explicitly trained with linguistic annotations grounded in such formalisms.

Much prior work has analyzed transformer representations using probing tasks based mainly on Universal Dependencies (UD) (Hewitt and Liang, 2019; Tenney et al., 2019). While valuable, the UD framework aims for cross-lingual consistency and potentially abstracts away from linguistic phenomena better captured by language-specific or theory-specific grammars. Indian languages (IL), characterized by rich morphology and relatively free word order, are often analyzed using frameworks derived from the classical Paninian grammatical tradition (Bharati et al., 1995). This tradition introduces distinctive linguistic concepts such as Kārakas (syntactico-semantic roles) and emphasizes morphological cues such as Vibhakti markers as fundamental in signaling syntactic structures.

A detailed overview of the specific Paninian annotation scheme (Begum et al., 2007) relevant to this work is provided in Appendix A.

In this work, we carry out the first Paninian analysis of multilingual transformer representations across layers. We probe the hidden layers of XLM-RoBERTa, IndicBERT-v2, and mBERT across seven Indian languages annotated under a shared Paninian annotation scheme (Begum et al., 2007). To measure these linguistic representations, we employ edge probing (Tenney et al., 2019) for syntactic dependency structures and token-classification probes for morphosyntactic features (Universal POS, Vibhakti markers).

Specifically, the paper aims to answer the following research questions clearly:

1. Do multilingual transformer models implicitly

124

represent Paninian dependency structures and morphosyntactic information?

2. If such representations exist, at which layer do they emerge distinctly?

3. How consistent are these observed patterns across different multilingual transformer architectures and different Indian languages?

Our key findings can be summarized as follows:

1. Paninian dependency structures are indeed implicitly encoded, predominantly emerging in middle layers of multilingual transformers (layers 6–9), in line with prior UD-based studies.

2. Lexical-level morphosyntactic signals (UPOS, Vibhakti markers) peak slightly earlier in the layers, highlighting differentiated storage of structural versus lexical linguistic information.

3. Despite shared general trends, we observe substantial variations across different architectures and languages, reflecting cross-linguistic diversity, data differences, and model-specific factors.

## 2 Experimental Setup

To investigate and compare the layer-wise encoding of Paninian grammar in different multilingual transformers, we conducted probing experiments using three base models: XLM-RoBERTa-base (Conneau et al., 2019), Multilingual BERT (Devlin et al., 2019) and IndicBERT-v2 (Doddapaneni et al., 2023). For each base model, parameters were kept frozen[1], and we extracted hidden states from all layers (including embeddings; Layer 0 up to Layer 12 for base models). Our analysis covers seven Indian languages: Hindi (hi), Kannada (ka), Malayalam (ml), Marathi (mr), Telugu (te), Urdu (ur), and Bengali (be). The datasets utilize annotations following an extended version of the Paninian dependency annotation scheme proposed by Begum et al. (2007).

**Probing Tasks and Models:** We designed diagnostic classifiers (Tenney et al., 2019) for three distinct tasks reflecting key aspects of the formalism. Following best practices for probing (Hewitt and Liang, 2019), we aimed for low-capacity probes:

[1]Only probe parameters were trained (Alain and Bengio, 2016; Hewitt and Liang, 2019).

- **Dependency Relations (Edge Probing):** To assess structural syntactic encoding, we trained probes using a biaffine attention/classifier architecture (Dozat and Manning, 2017). This probe incorporates non-linearity through intermediate single-layer MLPs (dimensionality 256, GELU activation) applied to the input hidden states before the biaffine transformations for predicting heads and relations.

- **UPOS Tags:** To capture basic lexical categories, we trained strictly linear classifier probes mapping the hidden state directly to tag logits.

- **Vibhakti Features:** Similarly, we used linear classifier probes to predict grammatical case/postposition markers. Vibhakti labels were extracted from the FEATS column based on observed patterns; non-applicable tokens were ignored during accuracy calculation.

Separate probe models were trained for each task, base model, language, and layer combination.

**Data Handling:** Input sequences were processed using the respective tokenizer for each base model (XLM-R, mBERT, IndicBERT-v2). Sequences exceeding a tokenized length of 128 tokens were filtered out. Remaining sequences were padded/truncated to 128 tokens. Target labels were aligned to the first sub-word token, with others ignored. Detailed dataset statistics are provided in Appendix B.

**Training and Evaluation:** Probes were trained independently for 3 epochs using the AdamW optimizer (Loshchilov and Hutter, 2017) with CrossEntropyLoss. Learning rates and batch sizes were tuned per model and task type based on preliminary experiments (Edge Probes: LR=1e-4, BS=64; Token Probes: LR=1.5e-3, BS=256). The model state yielding the best validation performance (LAS for dependencies, Accuracy for features) was selected. Evaluations were performed on held-out validation sets.

## 3 Results

We evaluated the performance of probes trained on each layer's representation across the three base models (XLM-RoBERTa, mBERT, IndicBERT-v2),

seven languages, and three probing tasks (Dependency Relations, UPOS, Vibhakti). Key findings are presented below.

**Dependency Relations (LAS):** Figure 1 presents the best Labeled Attachment Score (LAS) achieved by dependency probes trained on each layer (0-12) for the three base models across the seven Indian languages. A consistent trend across all three models and most languages is the distribution of syntactic information across layers. LAS is low at the initial embedding layer (L0), increases through the lower layers, generally peaks in the middle-to-upper-middle layers (typically Layers 7-9), and then declines towards the final layer (L12). This confirms the widely observed phenomenon that structural syntactic information is most saliently represented in the intermediate representations of transformer models (Jawahar et al., 2019; Tenney et al., 2019).

Model-wise, XLM-RoBERTa (Fig. 1a) and mBERT (Fig. 1b) achieve comparable peak performance on higher-resource languages like Hindi (peak LAS 61.6% for XLM-R at L8; 59.1% for mBERT at L7) and Urdu (peak LAS 48.1% for XLM-R; 47.9% for mBERT, both at L7). IndicBERT-v2 (Fig. 1c), despite its Indic pre-training, does not consistently outperform the general multilingual models on this task for all languages. While reaching a high peak for Hindi (62.1% at L6), its LAS scores for several other languages (ka, ml, be, mr, te) are often lower than those achieved by XLM-R or mBERT, particularly in the upper layers. IndicBERT-v2 shows earlier peaks in several languages (e.g., L6 for Hindi compared to L8 and L7 for XLM-R and mBERT respectively).

**UPOS Tagging (Accuracy):** Figure 2 illustrates the layer-wise accuracy for predicting Universal Part-of-Speech (UPOS) tags. Generally, UPOS tagging accuracy is high, significantly exceeding LAS scores, confirming that basic lexical category information is more readily extractable. Performance tends to peak relatively early compared to LAS, often plateauing across several lower-to-middle layers (e.g., Layers 3-9 for many languages in XLM-R and mBERT) before potentially declining slightly in the final layers.

Model-wise, XLM-RoBERTa and mBERT show strong results, with peak accuracies often reaching 85-90%+ for languages like Hindi and Kannada. IndicBERT-v2 generally achieves somewhat lower

peak accuracies on this task compared to the other two models (e.g., peaking around 75% for Hindi).

**Vibhakti Feature Prediction (Accuracy):** The accuracy for predicting Vibhakti features, a key morphological cue in Paninian grammar, is presented in Figure 3. Overall accuracy for this task is generally higher than LAS but can be lower than UPOS accuracy. Hindi and Urdu consistently yield the highest accuracies, frequently exceeding 80-90% across several layers in XLM-RoBERTa and mBERT, and maintaining high accuracy across almost all layers in IndicBERT-v2.

Compared to the sharper LAS peaks, strong Vibhakti prediction performance often extends across a broader range of middle and sometimes upper layers (e.g., Layers 2-10 for Hindi/Urdu in XLM-R/mBERT). The optimal layers for Vibhakti tend to overlap with or slightly precede the peak layers for dependency relations. Among the models, IndicBERT-v2 demonstrates particularly strong and stable Vibhakti prediction for Hindi and Urdu across nearly all layers.

**Cross-Lingual Variation:** Performance varies markedly across languages within each model (Figures 1, 2 and 3). While a comprehensive cross-lingual performance comparison is complicated by factors such as differing script representations and varying vocabulary coverage for each language within the pre-trained models, clear tiers of performance emerge. Hindi and Urdu consistently show the strongest LAS results, suggesting the models capture their Paninian structures relatively effectively. Kannada, Malayalam, Bengali, and Marathi form a mid-tier group, with peak LAS typically ranging from ~20% to ~40% depending on the model and language. Telugu consistently exhibits the lowest LAS scores across all models and layers (peak < 10% for XLM-R/mBERT, < 5% for IndicBERT-v2), a finding that strongly correlates with its significantly smaller probing dataset size (Appendix B) and potential underlying data sparsity in the models' pre-training.

## 4 Analysis and Discussion

Our multi-model, multi-task probing experiments reveal consistent layer-wise patterns for encoding Paninian grammar, alongside notable performance variations.

**Layer Specialization for Paninian Grammar:** Across all three models, a functional specialization

## Figure 1

**(a) XLM-RoBERTa**

| Layer | hi | ur | ka | ml | be | mr | te |
|---|---|---|---|---|---|---|---|
| 0 | 31.0 | 12.1 | 14.2 | 10.1 | 6.6 | 7.2 | 2.4 |
| 1 | 39.4 | 24.3 | 20.5 | 14.2 | 9.3 | 11.0 | 2.2 |
| 2 | 48.8 | 31.6 | 26.4 | 19.2 | 14.8 | 15.3 | 3.4 |
| 3 | 53.0 | 36.1 | 29.9 | 21.7 | 17.3 | 17.0 | 3.7 |
| 4 | 56.2 | 38.6 | 33.7 | 24.9 | 18.9 | 20.5 | 4.0 |
| 5 | 59.3 | 43.2 | 36.6 | 28.1 | 23.7 | 22.4 | 4.6 |
| 6 | 60.8 | 43.7 | 37.7 | 30.0 | 24.1 | 22.2 | 5.5 |
| 7 | 61.5 | 48.1 | 38.6 | 29.5 | 23.0 | 22.4 | 4.9 |
| 8 | 61.6 | 46.8 | 39.1 | 29.7 | 23.0 | 23.3 | 6.2 |
| 9 | 60.9 | 41.2 | 36.2 | 27.0 | 21.8 | 21.4 | 5.0 |
| 10 | 59.6 | 35.7 | 34.7 | 24.8 | 19.4 | 18.8 | 5.3 |
| 11 | 57.4 | 35.3 | 33.4 | 22.8 | 15.0 | 16.2 | 5.3 |
| 12 | 28.5 | 10.5 | 14.1 | 10.0 | 5.8 | 7.0 | 2.2 |

**(b) mBERT**

| Layer | hi | ur | ka | ml | be | mr | te |
|---|---|---|---|---|---|---|---|
| 0 | 35.1 | 21.1 | 17.3 | 12.6 | 8.3 | 11.1 | 4.3 |
| 1 | 39.4 | 22.5 | 20.4 | 15.0 | 10.1 | 10.7 | 3.9 |
| 2 | 45.0 | 28.2 | 22.1 | 15.8 | 12.3 | 13.0 | 4.1 |
| 3 | 51.0 | 36.2 | 29.6 | 22.0 | 16.4 | 18.5 | 4.7 |
| 4 | 54.2 | 41.0 | 33.1 | 25.1 | 19.3 | 22.0 | 4.3 |
| 5 | 55.7 | 43.9 | 35.2 | 27.2 | 20.2 | 24.3 | 5.0 |
| 6 | 58.0 | 46.3 | 37.7 | 30.3 | 23.1 | 27.4 | 5.3 |
| 7 | 59.1 | 47.9 | 38.0 | 31.7 | 24.5 | 27.3 | 5.9 |
| 8 | 58.6 | 47.1 | 37.7 | 31.2 | 23.4 | 26.9 | 5.8 |
| 9 | 57.6 | 45.7 | 37.0 | 30.5 | 22.7 | 25.7 | 4.7 |
| 10 | 56.2 | 42.4 | 35.1 | 29.1 | 21.9 | 22.8 | 5.6 |
| 11 | 54.3 | 36.9 | 32.7 | 26.9 | 18.3 | 21.0 | 4.9 |
| 12 | 46.6 | 25.8 | 25.4 | 19.9 | 12.4 | 13.4 | 2.7 |

**(c) IndicBERT-v2**

| Layer | hi | ur | ka | ml | be | mr | te |
|---|---|---|---|---|---|---|---|
| 0 | 43.0 | 17.5 | 25.5 | 17.3 | 11.7 | 11.3 | 2.8 |
| 1 | 52.4 | 22.5 | 30.6 | 20.6 | 14.4 | 12.9 | 3.7 |
| 2 | 57.6 | 28.8 | 34.1 | 24.6 | 16.7 | 15.9 | 3.4 |
| 3 | 59.7 | 31.5 | 36.9 | 26.3 | 17.9 | 15.6 | 3.6 |
| 4 | 61.2 | 36.3 | 38.6 | 28.8 | 20.1 | 17.6 | 3.3 |
| 5 | 62.1 | 36.6 | 39.5 | 30.2 | 20.3 | 18.8 | 4.0 |
| 6 | 62.1 | 36.4 | 40.0 | 30.9 | 20.8 | 19.1 | 3.7 |
| 7 | 61.0 | 34.4 | 39.5 | 29.7 | 20.3 | 17.9 | 3.6 |
| 8 | 59.7 | 31.1 | 38.0 | 28.4 | 20.5 | 17.4 | 3.4 |
| 9 | 57.7 | 29.2 | 36.2 | 27.7 | 18.9 | 17.1 | 3.1 |
| 10 | 55.3 | 26.4 | 35.0 | 25.5 | 17.5 | 16.8 | 3.7 |
| 11 | 32.8 | 12.8 | 21.3 | 15.8 | 7.5 | 6.9 | 2.5 |
| 12 | 33.2 | 12.9 | 18.8 | 13.3 | 7.0 | 7.0 | 1.9 |

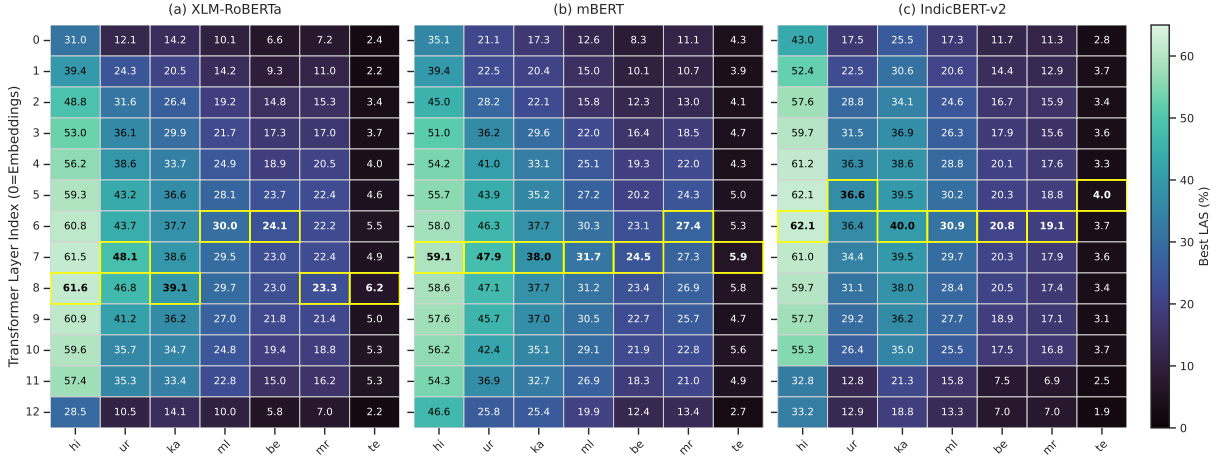*(Y axis: Transformer Layer Index (0=Embeddings); color bar: Best LAS (%))*

Figure 1: Layer-wise Dependency LAS (%) across Languages for (a) XLM-RoBERTa, (b) mBERT, and (c) IndicBERT-v2. Brighter colors indicate higher LAS. Yellow boxes highlight approximate peak performance regions for selected languages.

## Figure 2

**(a) XLM-Roberta**

| be | hi | ka | ml | mr | te | ur |
|---|---|---|---|---|---|---|
| 58.00 | 80.88 | 63.97 | 60.41 | 51.46 | 49.65 | 67.37 |
| 60.21 | 85.14 | 73.27 | 65.52 | 40.87 | 59.07 | 70.16 |
| 59.30 | 88.66 | 76.89 | 64.94 | 54.66 | 54.42 | 69.77 |
| 64.58 | 89.57 | 81.38 | 68.05 | 35.14 | 56.28 | 70.68 |
| 64.74 | 90.10 | 79.87 | 70.48 | 53.26 | 50.23 | 71.06 |
| 58.36 | 91.04 | 80.04 | 70.92 | 46.67 | 52.91 | 66.34 |
| 61.42 | 90.81 | 79.56 | 69.42 | 52.43 | 53.37 | 71.44 |
| 62.44 | 88.81 | 76.81 | 69.25 | 43.01 | 58.95 | 70.08 |
| 58.85 | 89.59 | 79.42 | 67.66 | 49.70 | 48.02 | 68.54 |
| 55.48 | 89.76 | 77.89 | 68.77 | 24.47 | 62.09 | 68.29 |
| 65.11 | 88.09 | 77.70 | 67.10 | 25.26 | 47.56 | 68.21 |
| 54.39 | 86.76 | 77.09 | 61.30 | 38.39 | 24.42 | 64.59 |
| 51.06 | 68.07 | 53.16 | 64.75 | 17.42 | 21.74 | 47.55 |

**(b) mBERT**

| be | hi | ka | ml | mr | te | ur |
|---|---|---|---|---|---|---|
| 60.99 | 72.30 | 73.44 | 67.38 | 53.51 | 55.93 | 71.87 |
| 61.07 | 71.65 | 74.38 | 67.99 | 64.83 | 59.07 | 72.44 |
| 61.74 | 75.19 | 75.62 | 67.24 | 65.12 | 55.70 | 73.78 |
| 64.00 | 73.11 | 82.36 | 68.97 | 67.73 | 52.44 | 76.75 |
| 64.69 | 75.50 | 84.44 | 72.22 | 65.80 | 65.70 | 77.32 |
| 64.97 | 75.10 | 85.57 | 71.89 | 56.78 | 65.35 | 77.89 |
| 61.62 | 73.00 | 85.40 | 73.12 | 51.46 | 57.67 | 76.97 |
| 63.07 | 66.37 | 85.65 | 72.47 | 65.65 | 67.09 | 73.24 |
| 62.14 | 71.11 | 85.24 | 73.04 | 65.41 | 67.91 | 72.74 |
| 63.05 | 68.99 | 85.59 | 73.13 | 67.82 | 58.72 | 73.19 |
| 61.37 | 68.04 | 85.10 | 73.02 | 61.32 | 61.98 | 72.66 |
| 59.39 | 67.23 | 82.31 | 69.60 | 64.57 | 58.14 | 69.72 |
| 48.10 | 57.16 | 75.24 | 61.02 | 51.41 | 45.93 | 61.36 |

**(c) IndicBERT-v2**

| be | hi | ka | ml | mr | te | ur |
|---|---|---|---|---|---|---|
| 57.41 | 68.30 | 67.96 | 61.43 | 39.32 | 52.79 | 32.77 |
| 57.37 | 73.49 | 68.62 | 61.27 | 50.30 | 53.37 | 33.96 |
| 56.71 | 74.16 | 68.22 | 65.10 | 27.55 | 53.37 | 34.75 |
| 53.17 | 73.12 | 68.62 | 64.74 | 40.08 | 59.07 | 34.41 |
| 56.34 | 73.57 | 67.93 | 63.51 | 43.63 | 58.84 | 34.26 |
| 59.88 | 73.64 | 66.71 | 61.81 | 48.88 | 60.47 | 33.74 |
| 50.69 | 74.11 | 67.68 | 63.27 | 28.12 | 55.70 | 32.96 |
| 58.85 | 73.31 | 67.26 | 62.87 | 28.10 | 53.95 | 34.57 |
| 59.07 | 73.55 | 66.26 | 63.75 | 48.66 | 61.05 | 33.50 |
| 55.22 | 73.62 | 65.63 | 63.26 | 34.46 | 55.93 | 33.82 |
| 58.69 | 75.42 | 69.78 | 63.73 | 57.37 | 54.88 | 32.37 |
| 58.71 | 74.50 | 67.45 | 61.29 | 30.62 | 48.14 | 33.42 |
| 54.27 | 72.74 | 65.56 | 61.46 | 22.61 | 53.37 | 34.10 |

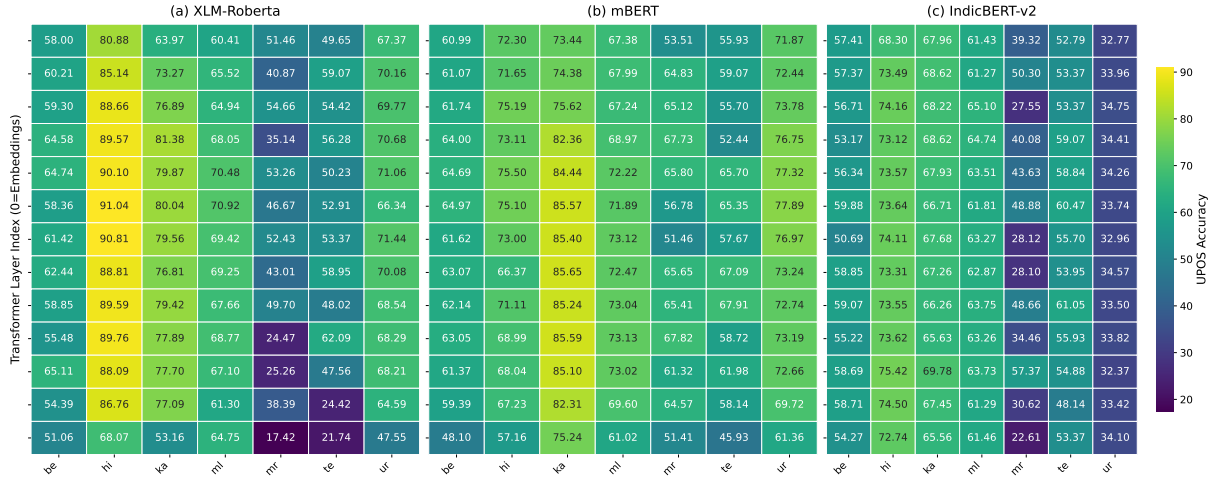*(Y axis: Transformer Layer Index (0=Embeddings); color bar: UPOS Accuracy)*

Figure 2: Layer-wise UPOS Accuracy (%) across Languages for (a) XLM-RoBERTa, (b) mBERT, and (c) IndicBERT-v2. Brighter colors indicate higher accuracy.

of layers aligns with prior probing studies on formalisms like Universal Dependencies (UD) (Jawahar et al., 2019; Tenney et al., 2019). Basic lexical information (UPOS tags) becomes accurately predictable in lower-to-middle layers. Morphological features (Vibhakti) also show strong representation across middle layers, their optimal encoding often overlapping with or slightly preceding layers most informative for syntactic structure. Crucially, complex Paninian dependency relations (LAS) consistently peak later, in the upper-middle layers (7-9), suggesting a hierarchical process where models integrate lexical/morphological cues to build syntactic representations. Performance generally degrades in final layers, possibly as representations specialize towards pre-training objectives.

**Model Architectures and Pre-training Influence:**
While layer-wise trends are broadly similar, absolute performance and peak locations vary across models. XLM-RoBERTa and mBERT show comparable LAS capabilities on higher-resource languages like Hindi and Urdu. IndicBERT-v2, despite its Indic-focused pre-training, does not uniformly outperform general multilingual models in LAS for all tested languages, though it achieves strong LAS for Hindi. However, IndicBERT-v2 excels in stable Vibhakti prediction for Hindi/Urdu across most layers, potentially reflecting better morphological encoding due to its specialized training. This nuanced behavior suggests language-family specific pre-training can enhance morphological representation, but generalization to complex syn-
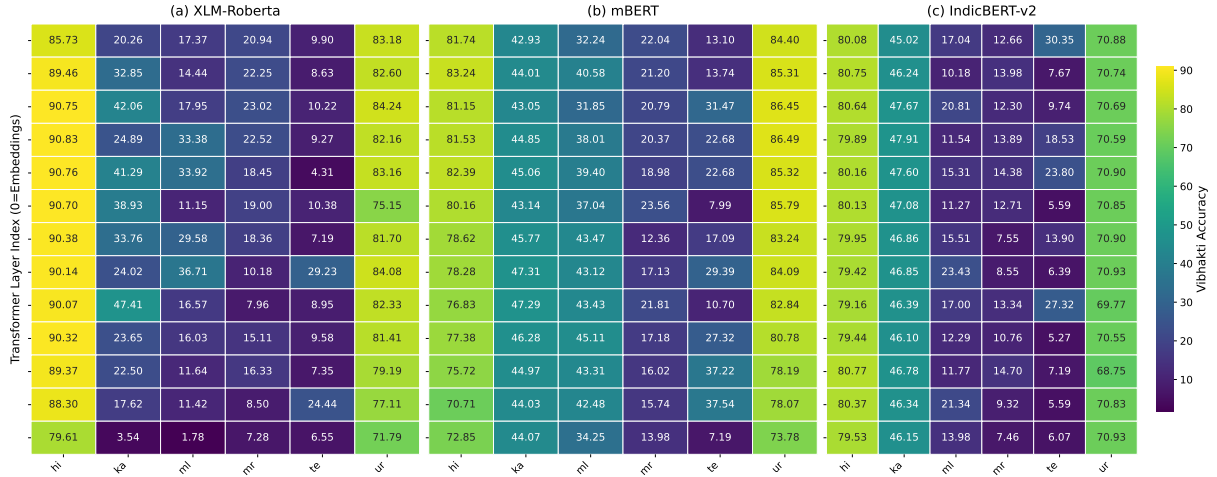
Figure 3: Layer-wise Vibhakti Accuracy (%) across Languages for (a) XLM-RoBERTa, (b) mBERT, and (c) IndicBERT-v2. Brighter colors indicate higher accuracy.

tax across diverse languages within that family remains challenging and interacts with other model properties.

**Cross-Lingual Consistency and Variation:** Significant cross-lingual variation in probing performance is evident for all tasks. While direct comparison is complicated by differing script representations and vocabulary coverage, Hindi and Urdu consistently yield the strongest results. The markedly lower performance for Telugu correlates with its smaller probing dataset size (Appendix B) and likely reflects underlying challenges from pre-training data sparsity or dataset quality. This underscores that probing performance reflects an interplay between information encoded by the base model and the characteristics of the probe training dataset.

**Encoding Paninian-Specific Information:** Our results demonstrate that diagnostic probes can successfully extract information pertinent to the Paninian grammatical framework – Kāraka-based dependency relations and Vibhakti features – from the frozen representations of these multilingual transformers. The ability to predict these suggests models implicitly learn representations sensitive to this distinct formalism, primarily consolidating structural knowledge in their middle layers.

## 5  Conclusion

We presented a layer-wise probing analysis comparing the encoding of Paninian grammatical information within XLM-RoBERTa, mBERT, and IndicBERT-v2 across seven Indian languages, examining dependency relations, UPOS tags, and Vibhakti features. Our findings reveal that Paninian dependency structure generally peaks in the upper-middle transformer layers, following lexical and morphological feature encoding in lower-to-middle layers, consistent with known patterns of linguistic representation.

Substantial cross-lingual and cross-model variations were observed. While IndicBERT-v2 showed strengths in Vibhakti prediction for core Indic languages, it did not uniformly surpass general multilingual models in representing Paninian dependency structures. Performance differences across languages correlate strongly with probing dataset sizes and likely reflect variations in pre-training data. Our results confirm that probing effectively reveals how theory-specific grammatical formalisms are represented within standard multilingual models.

## 6  Limitations

This study is limited to three base models and seven Indian languages; findings may not generalize broadly. Probe performance reflects both model representations and probe/dataset characteristics, with notable dataset size/quality variations (e.g., Telugu, Appendix B). The Paninian features probed (deprel, Vibhakti) are not exhaustive. Our analysis of frozen representations establishes correlations, not causal model mechanisms. Future work involves expanding model/language scope, probing more fine-grained Paninian features and studying the emergence of these representations.

# References

Guillaume Alain and Yoshua Bengio. 2016. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*.

Rafiya Begum, Samar Husain, Arun Dhwaj, Dipti Misra Sharma, Lakshmi Bai, and Rajeev Sangal. 2007. Dependency annotation scheme for indian languages. In *Proceedings of the Linguistic Annotation Workshop (LAW '07)*, pages 149–156, Prague, Czech Republic. Association for Computational Linguistics.

Yonatan Belinkov and James Glass. 2019. Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics (TACL)*, 7:49–72.

Akshar Bharati, Vineet Chaitanya, and Rajeev Sangal. 1995. *Natural Language Processing: A Paninian Perspective*. Prentice-Hall of India, New Delhi.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Sumanth Doddapaneni, Rahul Aralikatte, Gowtham Ramesh, Shreya Goyal, Mitesh M. Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2023. Towards leaving no indic language behind: Building monolingual corpora, benchmark and models for indic languages. *Preprint*, arXiv:2212.05409.

Timothy Dozat and Christopher D Manning. 2017. Deep biaffine attention for neural dependency parsing. In *International Conference on Learning Representations (ICLR)*.

John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.

Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does BERT learn about the structure of language? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.

# A Background

Our probing analysis leverages the Paninian dependency annotation scheme, specifically developed for accurately representing linguistic phenomena within morphologically rich and relatively free word order ILs (Begum et al., 2007).

The Paninian grammatical tradition, originating from the ancient linguist Pāṇini, describes sentences primarily through modifier-modified dependency structures centered around the verb (Bharati et al., 1995). A crucial aspect of Paninian grammar is the use of Kārakas, specialized syntactico-semantic relations linking verbs and their arguments or modifiers. In the annotation scheme adopted here (Begum et al., 2007), six main Kārakas are identified:

1. *adhikaraṇa* (location)

2. *apādān* (source)

3. *sampradān* (recipient)

4. *karaṇa* (instrument)

5. *karma* (theme, loosely object-like)

6. *karta* (agent, loosely subject-like)

Importantly, Kārakas are not exact equivalents to purely semantic thematic roles (e.g., agent, patient). Instead, they encode a distinctly Paninian syntactico-semantic perspective. Consider for instance the English-like example 'key opened the door'. While semantically an instrument, the 'key' here would be annotated as the *karta* (loosely agent-like) Kāraka in the Paninian tradition (Bharati et al., 1995).

Identification of these Kārakas depends heavily on morphological indicators such as Vibhakti (case-endings, postpositions) and verb-based TAM markers (tense-aspect-modality) within sentences (Begum et al., 2007). The strong, systematic correlation between morphological features and syntactic dependency structures motivates our probing

approach: we probe multilingual transformer models' representations for both the structural Kāraka relations (via dependency links and edge labels) and essential morphological cues (*Vibhaktis*, UPOS tags), analyzing explicitly how these linguistic representations are distributed layer-wise.

## A.1 Vibhakti

In the context of Paninian grammar and many modern Indian languages, *Vibhakti* refers to morphological markers, primarily case endings or postpositions, that are attached to nouns or noun phrases. These markers play a crucial role in signaling the grammatical function and syntactico-semantic role of the noun phrase within the sentence.

While often translated loosely as *case*, Vibhakti in the Paninian tradition is intimately linked to the concept of *Kārakas* (described in Appendix A). Specific Vibhaktis are typically associated with signaling specific Kāraka roles (e.g., a particular Vibhakti might commonly mark the *karta* 'agent-like' role, while another marks the *karma* 'theme-like' role, and others mark instrument, location, etc.). However, the mapping is not always one-to-one and can be influenced by other factors like verb semantics and sentence structure.

Essentially, Vibhaktis provide explicit surface cues about the underlying grammatical relationships in the sentence, making them particularly important in languages with relatively flexible word order where syntactic function is not solely determined by position.

## B Dataset Statistics

Table 1 provides statistics for the annotated datasets used in our probing experiments. Counts reflect the data *after* filtering sequences longer than 128 tokens but *before* any potential subsetting for development runs. Please note that the Bengali annotated data does not include Vibhakti features.

| Language | | Sentences | | Tokens | |
|---|---|---|---|---|---|
| Code | Name | train | val | train | val |
| be | Bengali | 3939 | 488 | 99504 | 12389 |
| hi | Hindi | 19855 | 2478 | 570772 | 71333 |
| ka | Kannada | 10388 | 1297 | 263827 | 32910 |
| ml | Malayalam | 7109 | 882 | 178754 | 21873 |
| mr | Marathi | 3608 | 451 | 91077 | 11461 |
| te | Telugu | 1514 | 185 | 19663 | 2219 |
| ur | Urdu | 4871 | 607 | 164129 | 21606 |

Table 1: Statistics of the Paninian–annotated datasets used for probing (post-filtering, pre-subsetting).