

Principal Parts Detection for Computational Morphology: Task, Models and Benchmark

Anonymous ACL submission

Abstract

Principal parts of an inflectional paradigm, defined as the minimal set of paradigm cells required to deduce all others, constitute an important concept in theoretical morphology. This concept, which outlines the minimal memorization needed for a perfect inflector, has been largely overlooked in computational morphology despite impressive advances in the field over the last decade. In this work, we posit PRINCIPAL PARTS DETECTION as a computational task and construct a multilingual dataset of verbal principal parts covering ten languages, based on Wiktionary entries. We evaluate an array of PRINCIPAL PARTS DETECTION methods, all of which follow the same schema: characterize the relationships between each pair of inflectional categories, cluster the resulting vector representations, and select a representative of each cluster as a predicted principal part. Our best-performing model, based on Edit Script between inflections and using Hierarchical K-Means, achieves an F1 score of 55.05%, significantly outperforming a Random Baseline of 21.20%. While our results demonstrate that some success is achievable, further work is needed to thoroughly solve PRINCIPAL PARTS DETECTION, a task that may be used to further optimize inputs for morphological inflection, and to promote research into the theoretical and practical importance of a compact representation of morphological paradigms.

1 Introduction

Morphological analysis is essential for understanding natural language, particularly in languages with complex inflectional systems. In both linguistic theory and language pedagogy, the concept of *principal parts* plays a central role in structuring and simplifying inflectional paradigms (Finkel and Stump, 2007; Stump and Finkel, 2013). Principal parts form the minimal subset of paradigm cells from which all other forms can be systematically derived.

By identifying these key forms, principal parts provide a compact representation of inflection tables and facilitate the analysis of morphologically rich languages. Despite their theoretical significance, the detection of principal parts remains largely unexplored in computational morphology. While they have inspired research in inflection and reinflection (Cotterell et al., 2017; Liu and Hulden, 2020), they are rarely used explicitly. Most computational approaches instead rely on a single citation form, the lemma (Cotterell et al., 2016; Goldman et al., 2023), or select input forms randomly (Cotterell et al., 2016; Kann et al., 2017). This reliance on suboptimal input representations overlooks the potential of principal parts as a more efficient foundation for inflectional modeling.

In this work, we introduce PRINCIPAL PARTS DETECTION as a formal task within computational morphology. Given a large collection of inflection tables, the goal is to determine which paradigm cells constitute the minimal principal-part set. Crucially, inflection tables typically contain standard morphological annotations but are not explicitly labeled with principal parts, making this an unsupervised learning problem. To promote research in this area, we deliver a standardized dataset covering the verbal paradigms of ten diverse languages. We sourced Principal parts for each language from online dictionaries, where they are often listed to aid language learners, and obtained full inflection tables from UniMorph (Batsuren et al., 2022).

We develop several computational approaches for PRINCIPAL PARTS DETECTION, leveraging the defining property of principal parts: their predictable and systematic relationships with other forms in the paradigm. Our models characterize inter-cell similarity and cluster inflected forms into *sub-paradigms*, selecting a representative cell from each sub-paradigm as candidate principal parts. We explore different methods for *characterizing* inter-cell relations, including Edit Distance, Edit

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

Script, and Reinflection Accuracy, and we experiment with *clustering* techniques such as Affinity Propagation and a modified K-Means algorithm. Our best-performing system, using Edit Script similarity measure + Hierarchical K-Means clustering, achieves an average F1 score of 55.05% across the ten languages in our dataset, significantly outperforming a Random Baseline of 21.20%.

By formalizing PRINCIPAL PARTS DETECTION as a computational task, we lay the groundwork for future research on more efficient morphological representations. To the best of our knowledge, this is the first work to deliver a standardized benchmark of PRINCIPAL PARTS DETECTION alongside a fully-operational detection framework. Successfully solving this task could enhance applications in morphological inflection and analysis by providing more informative input forms. Our findings suggest that principal parts can be computationally identified with reasonable accuracy, but further improvements are necessary to fully realize their potential.

2 The PRINCIPAL PARTS DETECTION Task and Dataset

The PRINCIPAL PARTS DETECTION Task. The task of PRINCIPAL PARTS DETECTION is defined as identifying the minimal set of cells within a paradigm that, when known, allow the derivation of all other paradigm forms. For instance, in English, the principal parts of the verbal paradigm are the cells corresponding to the infinitive, simple past and past participle (for example, *eat*, *ate*, and *eaten*), as these forms are not predictable from one another, especially for strong verbs. On the other hand, the forms corresponding to the present participle and the 3rd person singular present are deterministically predictable from the infinitive and they therefore provide no additional information for inflection if the infinitive is known.

Formally, the task of PRINCIPAL PARTS DETECTION is defined given a language L , a paradigm $P \in L$, and a large set of inflection tables $T = \{t_{P,1}^L, t_{P,2}^L, \dots, t_{P,n}^L\}$ that belong to that paradigm. The goal is to identify a minimal set of cells $C_{PP} \subseteq P$ from which all other cells in all inflection tables of P can be accurately deduced.

The PRINCIPAL PARTS DETECTION Dataset. In order to empirically assess methods for the detection of principal parts, we first need to have a dataset to evaluate against. To this end, we

constructed the PRINCIPAL PARTS DETECTION dataset, containing ten typologically diverse languages, where every paradigm in every language is characterized by a set of target principal parts that systems can be evaluated against.

The input side of the task contains complete inflection tables, based on the UniMorph corpus (Batsuren et al., 2022), which provides inflection tables for 168 languages organized by lexeme and morpho-syntactic features. We sourced gold principal parts — that are the desired output — from a combination of Wiktionary and other trusted online dictionaries or language teaching websites. Based on the availability of data sources for both input and output, we selected ten typologically diverse languages: Hebrew, English, French, German, Spanish, Danish, Swedish, Finnish, Turkish and Latin.

The dataset preparation process involved normalizing the data for consistency across languages. Redundant and derivational forms were excluded, leaving only core inflectional forms. Inconsistent feature sets were removed, and problematic entries from the original sources were manually corrected to ensure a reliable dataset (for more details, see Appendix A). The PRINCIPAL PARTS DETECTION dataset provides a strong foundation for computational models, bridging linguistic theory and practical applications. By curating this multilingual dataset, we ensure a robust resource for future research in morphological inflection. The next section shifts focus to computational methods for detecting principal parts, drawing on the linguistic insights outlined in the literature.¹

3 Translating Linguistic Insights into Computational Methods

The core linguistic principle underlying PRINCIPAL PARTS DETECTION is that principal parts encapsulate the implicative relationships that exist between cells in inflectional paradigms, allowing a small set of cells to reconstruct the full inflectional table. To translate this principle into a computationally tractable problem, we frame PRINCIPAL PARTS DETECTION as the automatic identification of a minimal, generative subset of paradigm cells that can generate all other cells via these implicative relationships.

We hypothesize that cells of different feature sets in an inflection table exhibit measurable simi-

¹The data is publicly available in <https://www.will.be.released/upon.acceptance>.

larities in their realized surface forms, that in turn reflect morphological and structural relationships between these feature-set cells. By capturing patterns of inter-dependence between cells, we approximate the implicative structure of a paradigm *without* relying on explicit linguistic annotations of principal parts. To systematically model these inter-dependencies, we introduce the notion of a *sub-paradigm*, which is essentially a sub-set of paradigm cells with shared implicative properties, and define distinct areas-of-interdependence within a paradigm. Although sub-paradigms are not a formal linguistic concept, they provide a structured way to model implicative relationships computationally, facilitating the detection of principal parts.

This conceptualization leads to a three-phased methodology for PRINCIPAL PARTS DETECTION. First, we *characterize* the relationships between pairs of cells by computing similarity measures that capture their surface and structural dependencies. Next, we *cluster* inter-related cells into sub-paradigms, each of which will be represented by a single principal part in the final set. Finally, we *select* one representative feature set per sub-paradigm as its designated principal part, ensuring maximal coverage of the paradigm cells with minimal redundancy. The instantiation of these (i) *characterization*, (ii) *clustering* and (iii) *candidate selection* phases gives rise to a host of PRINCIPAL PARTS DETECTION concrete implementations that we can define and empirically assess — as we discuss next.

4 Framework and Task Empirical Design

The PRINCIPAL PARTS DETECTION framework we propose here is composed of three interconnected stages: characterization, clustering, and principal parts selection, each implemented using well-defined computational methods. These stages operate independently, meaning that different configurations of the framework can mix and match methods in seeking the best combination. Let us briefly review the computational models we consider for the different phases.

4.1 Characterization: Quantifying Relationships Between Feature Sets

The characterization stage quantifies the relationships between paradigm cells by computing numerical similarity scores between them. This work explores three distinct characterization methods, offering a different perspectives on the relation be-

tween cells.

Edit Distance A metric that measures surface-level similarity between forms based on minimal edit operations — insertions, deletions, or substitutions — required to transform one form into another (Levenshtein, 1966). This method is implemented by computing the average Edit Distance from each feature set to all others (calculated across all their surface realizations), treating one as the source and the rest as destinations. The resulting vector representations store these averaged distances, capturing the surface-level similarity between feature sets. Pairs of paradigm cells with low Edit Distance scores exhibit orthographic overlap.

Edit Script A metric that captures transformational diversity by analyzing character-level transformations between paradigm cells. Unlike traditional Edit Script approaches (Wagner and Fischer, 1974; Myers, 1986), which focus on the exact sequence of operations needed to transform one string into another, this approach computes the number of unique character-level transformations observed across all surface realizations of each feature set pair. Each transformation is counted only once per feature set pair (calculated across all their surface realizations), capturing distinct transformational patterns rather than repeatedly occurring character changes. The result is a vector representation for each feature set pair, where each entry encodes the number of unique transformations required to convert one feature set to another, representing their transformational distance. This method provides insight into the variation in morphological transformations within a paradigm. Feature sets with lower transformation diversity may exhibit more stable morphological patterns, making them stronger principal part candidates. In contrast, higher transformation diversity may signal greater variability in inflectional behavior, which can affect predictability within the paradigm.

Reinflection Accuracy A metric that evaluates the functional predictability of feature sets. It leverages a neural reinflection model trained to generate a target form given a source form and the morpho-syntactic features of the target. Unlike edit-based methods that focus on surface similarity and transformational diversity, Reinflection Accuracy captures the functional dependencies between feature sets, reflecting their predictive capacity within a

280	paradigm.	
281	Reinflection accuracy is particularly effective	330
282	in languages with complex inflectional systems,	331
283	where orthographic similarity alone is not a re-	332
284	liable predictor of implicative relationships. By	333
285	capturing functional dependencies rather than sur-	334
286	face transformations, it provides a direct measure	
287	of a feature set’s ability to generate other forms.	
288	However, its performance depends on training data	
289	quality and resource availability. In low-resource	
290	settings, data sparsity may lead to biased results,	
291	and the approach is computationally intensive, as it	
292	requires training multiple models—one model per	
293	feature set. Despite these challenges, its ability to	
294	model functional predictability makes it a valuable	
295	tool for identifying feature sets that serve as princi-	
296	pal parts, particularly in morphologically complex	
297	languages.	
298	Each characterization method produces a simi-	
299	larity table, where rows represent source feature	
300	sets and columns represent target feature sets, en-	
301	coding pairwise relationships (see Appendix B).	
302	Before clustering, all similarity matrices are stan-	
303	dardized by removing the mean and scaling to unit	
304	variance to ensure comparability across methods.	
305	These standardized characterization tables form the	
306	foundation for the clustering stage.	
307	4.2 Clustering: Structuring Feature Sets into	
308	Sub-Paradigms	
309	The clustering stage groups feature sets based on	
310	their quantified relationships, approximating sub-	
311	paradigms that reflect the internal organization of	
312	inflectional paradigms. The framework implements	
313	two clustering algorithms, each offering different	
314	advantages. As with characterization, only one	
315	clustering algorithm is used at a time.	
316	Affinity Propagation A message-passing clus-	
317	tering algorithm that dynamically determines the	
318	number of clusters based on pairwise similarity	
319	scores (Frey and Dueck, 2007). Unlike traditional	
320	clustering methods, it does not require a predefined	
321	number of clusters. Instead, it iteratively updates	
322	responsibility and availability values, which deter-	
323	mine how well a feature set serves as an exemplar	
324	(cluster center), until the algorithm converges on	
325	a final set of exemplars. This property makes it	
326	particularly well-suited for paradigms with high	
327	morphological variability. The algorithm is im-	
328	plemented using scikit-learn’s <code>AffinityPropagation</code>	
329	module, with similarity scores computed as nega-	
	tive squared Euclidean distances. The preference	330
	parameter is set to the median similarity value, al-	331
	lowing clusters to emerge naturally. Additional	332
	parameters include a convergence iteration limit of	333
	30 and a random state value of 10.	334
	Hierarchical K-Means A hierarchical variant	335
	of K-Means that recursively partitions feature sets	336
	into two clusters per iteration until a well-defined	337
	clustering structure is reached. The stopping cri-	338
	terion is determined using the Calinski–Harabasz	339
	Index (CHI) (Caliński and Harabasz, 1974), which	340
	evaluates clustering quality by comparing between-	341
	cluster dispersion to within-cluster cohesion. At	342
	each step, the CHI is computed across the entire	343
	clustering structure to evaluate how well-separated	344
	the clusters are relative to their internal cohesion.	345
	To prevent over-segmentation, clustering stops if	346
	the number of clusters in the new best CHI solu-	347
	tion exceeds that of the previous best CHI solu-	348
	tion by more than one cluster. The algorithm is im-	349
	plemented using scikit-learn’s <code>KMeans</code> module with a	350
	random state value of 10.	351
	By grouping feature sets into paradigm subsets,	352
	the clustering stage provides a data-driven approx-	353
	imation of sub-paradigms. The resulting clusters	354
	serve as inputs for the principal parts selection	355
	stage.	356
	4.3 Principal Parts Selection: Identifying	357
	Representative Feature Sets	358
	The principal parts selection stage finalizes the	359
	PRINCIPAL PARTS DETECTION framework by	360
	transforming clusters into a compact and generative	361
	summary of the paradigm. This stage selects one	362
	representative feature set per cluster, encapsulating	363
	its defining structural and transformational relation-	364
	ships. These feature sets collectively constitute the	365
	principal parts, providing comprehensive coverage	366
	while maintaining a balance between compactness	367
	and predictive capacity.	368
	Concretely, we use the Minimum Average In-	369
	flexional Length criterion. That is, the feature set	370
	with the minimal average inflectional length in its	371
	cluster, calculated across all its surface realizations,	372
	is chosen as the principal part. This ensures that	373
	the selected feature set is both efficient and central	374
	within its cluster. This selection criterion aligns	375
	with a linguistic insight that shorter inflectional	376
	paths often correspond to forms that are central	377
	within the paradigm, making them structurally sig-	378
	nificant within inflectional systems.	379

5 Experimental Setup and Results

We conduct a series of experiments to evaluate the effectiveness of the PRINCIPAL PARTS DETECTION framework across ten typologically diverse languages. The evaluation compares six model configurations, each formed by pairing one of three characterization methods—Edit Distance, Edit Script, and Reinflection Accuracy—with one of two clustering algorithms—Affinity Propagation and Hierarchical K-Means. To establish a performance threshold, we include a Random Baseline, which selects principal parts at random.

5.1 Dataset

The experiments are conducted on the PRINCIPAL PARTS DETECTION dataset, which comprises ten typologically diverse languages, divided into a development set (Hebrew, English, French, German, and Spanish) and a test set (Danish, Swedish, Finnish, Turkish, and Latin).

The development set represents varied morphological structures. Hebrew exhibits synthetic morphology, encoding multiple grammatical elements within single word forms. English, in contrast, is analytic, primarily relying on word order and auxiliary constructions for grammatical relationships. French and Spanish, as fusional languages, encode tense, mood, and person within single inflections, albeit with varying degrees of regularity. German, a hybrid case, incorporates both fusional and analytic morphological characteristics, presenting distinct patterns for analysis. This linguistic diversity ensures that models are trained on paradigms with different degrees of morphological richness, regularity, and complexity.

The test set is designed to assess generalization across languages with distinct inflectional systems. Finnish and Turkish exemplify agglutinative morphology, where grammatical meaning is expressed through concatenative morphemes. Latin, a highly inflected classical language, provides a rigorous test case for evaluating the models’ ability to handle case, number, and gender distinctions. Danish and Swedish, characterized by relatively regular inflectional systems, contribute typological variety while testing the models’ robustness in less complex paradigms.

By structuring the dataset to reflect a wide range of linguistic variation, this division ensures a comprehensive evaluation of the framework’s adaptability to diverse morphological systems and its

ability to generalize across typologically distinct languages.

5.2 Evaluation Metric

To evaluate model effectiveness, we use the F1 score, which balances precision and recall to assess both accuracy and completeness in PRINCIPAL PARTS DETECTION.

In addition to reporting F1 scores, we compare model performance against a Random Baseline, which selects principal parts randomly within each paradigm. Given a paradigm with x feature sets and y gold principal parts, the probability of randomly selecting a correct principal part is $\frac{y}{x}$. Since the baseline selects y principal parts, the expected number of correct predictions is $y \times \frac{y}{x} = \frac{y^2}{x}$. From this, the expected precision, recall, and F1 score are all: $F1 = \frac{y}{x}$. Since principal parts are inherently sparse within most paradigms, the Random Baseline represents a challenging threshold. Models that significantly exceed this score demonstrate an ability to detect principal parts systematically rather than relying on chance.

5.3 Reinflection Settings

For models utilizing Reinflection Accuracy, we train a separate reinflection model for each feature set, treating it as the source while all other feature sets serve as targets. The model is based on the Base LSTM architecture (Goldman et al., 2021), a character-based sequence-to-sequence model comprising a one-layer bidirectional LSTM encoder and a one-layer unidirectional LSTM decoder with a global soft attention layer (Bahdanau et al., 2014). Each model is trained for 50 epochs, optimizing categorical cross-entropy.

The dataset is split 70%-30%, ensuring that test lexemes are unseen during training. Each feature set is trained using a dedicated dataset, where it serves as the source inflection across different lexemes. Since each feature set is evaluated on its ability to generate all other feature sets within the paradigm, corresponding test sets are created—one per target feature set.

Each trained model is evaluated on how accurately it inflects from its assigned source feature set to each target feature set. The resulting accuracy scores form representation vector, capturing a feature set’s proficiency in generating others. Feature sets with high Reinflection Accuracy scores demonstrate strong predictive capacity, making them effective candidates for principal parts.

Model	Algorithmic Evaluation
Random Baseline	21.20
Edit Distance + Affinity Propagation	31.29
Edit Distance + Hierarchical K-Means	32.51
Reinlection Accuracy + Hierarchical K-Means	42.43
Edit Script + Affinity Propagation	44.62
Reinlection Accuracy + Affinity Propagation	45.56
Edit Script + Hierarchical K-Means	55.05

Table 1: Averaged F1 scores of PRINCIPAL PARTS DETECTION models across ten languages. The table compares different model configurations, highlighting the best-performing model.

5.4 Results

Table 1 presents the average F1 scores across ten languages, providing a comparative evaluation of model performance. All models outperform the Random Baseline, which achieves the lowest F1 score of 21.20%. The best-performing model, Edit Script + Hierarchical K-Means, achieves an F1 score of 55.05%, demonstrating its ability to effectively capture and cluster morphological patterns across diverse languages.

Reinlection Accuracy models perform competitively, with F1 scores of 45.56% (Affinity Propagation) and 42.43% (Hierarchical K-Means). In contrast, Edit Distance-based models yield lower scores of 31.29% and 32.51%, indicating that surface-level similarity alone is insufficient for PRINCIPAL PARTS DETECTION.

Overall, all tested methods surpass the Random Baseline by at least 10.09 points, with the best-performing model exceeding it by 33.85 points. These results confirm the effectiveness of the proposed methodology, demonstrating a substantial improvement over random selection.

Table 2 provides a language-specific breakdown of F1 scores, offering further insight into model performance across different morphological typologies. Edit Script + Hierarchical K-Means achieves top performance in Hebrew, French, Spanish, Turkish, and Latin, confirming its adaptability across different morphological systems. Reinlection Accuracy-based models perform particularly well in English, Spanish, Finnish, and Swedish, suggesting that functional predictability is well-suited for these languages.

Interestingly, while Reinlection Accuracy + Affinity Propagation ranks second overall (45.56%), it does not consistently outperform other models across all languages. In Danish and

Model	Hebrew	English	French	German	Spanish	Danish	Swedish	Finnish	Turkish	Latin
Random Baseline	20.68	60.00	14.28	16.66	2.53	62.50	26.30	2.48	28.00	6.25
Edit Distance + Affinity Propagation	33.30	66.70	37.50	46.20	15.40	57.10	40.00	0.00	0.00	16.70
Edit Distance + Hierarchical K-Means	25.00	57.10	44.40	44.40	0.00	57.10	57.10	0.00	0.00	40.00
Reinlection Accuracy + Hierarchical K-Means	25.00	85.70	44.40	28.60	50.00	57.10	43.50	50.00	0.00	40.00
Edit Script + Affinity Propagation	50.00	80.00	54.50	66.70	36.40	50.00	60.00	23.50	6.90	18.20
Reinlection Accuracy + Affinity Propagation	36.40	80.00	26.70	60.00	16.70	75.00	75.00	46.20	17.40	22.20
Edit Script + Hierarchical K-Means	50.00	80.00	54.50	60.00	50.00	72.70	60.00	33.30	50.00	40.00

Table 2: Language-specific F1 scores of PRINCIPAL PARTS DETECTION models across ten languages. The table highlights variations in model effectiveness across different morphological typologies.

Swedish, its relatively strong results suggest an advantage in regular inflectional systems where paradigmatic structures are highly predictable. Conversely, in fusional languages like Spanish, where single inflections encode multiple grammatical features, it faces challenges in PRINCIPAL PARTS DETECTION.

In contrast, Edit Distance-based models fail to rank highest in any language, reinforcing the conclusion that surface-level similarity alone is insufficient for PRINCIPAL PARTS DETECTION. These findings emphasize the importance of selecting appropriate characterization methods based on linguistic properties and show that transformational diversity (Edit Script) and functional predictability (Reinlection Accuracy) are particularly effective strategies.

6 Analysis

We analyze how methodological factors shape model performance, focusing on transformations in characterization data and the effectiveness of clustering strategies. This evaluation highlights structural patterns influencing clustering quality and examines the extent to which clustering results align with ideal principal parts selection.

6.1 Transpose Ablation: Evaluating the Impact of Data Orientation

The Transpose Ablation study investigates whether swapping the rows and columns of the characterization tables influences clustering quality and principal parts selection. This transformation is particularly relevant for Reinlection Accuracy, where the original tables encode directional relationships—rows indicate how easily a feature set can inflect from itself to others, while columns represent the reverse relationship. By transposing these tables, we test whether an alternative structural alignment improves performance.

Model	Transpose	Algorithmic Evaluation
Reinlection Accuracy + Affinity Propagation	✗	45.56
	✓	44.05
Reinlection Accuracy + Hierarchical K-Means	✗	42.43
	✓	43.14

Table 3: Algorithmic evaluation of Reinlection Accuracy models with and without transposition across ten languages. The table presents the averaged F1 scores for models before and after transposition, highlighting its varying impact depending on the clustering algorithm.

Transposition is applied only to Reinlection Accuracy models, as Edit Distance and Edit Script methods generate symmetric similarity matrices, making transposition redundant. We evaluate two models: Reinlection Accuracy + Affinity Propagation and Reinlection Accuracy + Hierarchical K-Means, comparing their performance before and after transposition.

The results in Table 3 show that transposition affects models differently. Reinlection Accuracy + Affinity Propagation experiences a slight decrease in performance (45.56% \rightarrow 44.05%), while Reinlection Accuracy + Hierarchical K-Means improves marginally (42.43% \rightarrow 43.14%). This suggests that transposition does not universally enhance clustering effectiveness and that its impact depends on the underlying clustering strategy.

Despite the minor improvement in Hierarchical K-Means, transposed results are excluded from the main evaluation due to their limited effect and misalignment with the principal parts definition. Since original (non-transposed) feature sets encode generative properties crucial for inflection, preserving this structure remains preferable. These findings suggest that alternative data transformations, better aligned with the linguistic task, may offer greater benefits.

6.2 Oracle Evaluation

To assess the theoretical upper limit of model performance, we conduct an Oracle evaluation, where principal parts are selected with perfect knowledge rather than through clustering. This evaluation distinguishes clustering effectiveness from principal parts selection quality, highlighting areas for improvement.

Table 4 reveals substantial gaps between Oracle and Algorithmic scores, underscoring clustering limitations and principal parts selection inefficiencies. Edit Script + Hierarchical K-Means achieves the highest Oracle score (76.21%), con-

Model	Evaluation	
	Oracle	Algorithmic
Edit Distance + Affinity Propagation	40.08	31.29
Edit Distance + Hierarchical K-Means	50.57	32.51
Reinlection Accuracy + Affinity Propagation	58.78	45.56
Reinlection Accuracy + Hierarchical K-Means	65.64	42.43
Edit Script + Affinity Propagation	54.16	44.62
Edit Script + Hierarchical K-Means	76.21	55.05

Table 4: Oracle and Algorithmic evaluations of PRINCIPAL PARTS DETECTION models across languages. Oracle evaluation assumes perfect knowledge of principal parts, establishing an upper bound on performance, while Algorithmic evaluation reflects actual model performance.

Model	Transpose	Evaluation	
		Oracle	Algorithmic
Reinlection Accuracy + Affinity Propagation	✗	58.78	45.56
	✓	58.51	44.05
Reinlection Accuracy + Hierarchical K-Means	✗	65.64	42.43
	✓	67.70	43.14

Table 5: Oracle and Algorithmic evaluations of Reinlection Accuracy before and after transposition. The table examines how transposition affects clustering quality under both ideal (Oracle) and algorithmic conditions.

firming strong clustering performance. However, the 21.16-point gap suggests that principal parts selection remains a limiting factor.

Conversely, Edit Distance + Affinity Propagation exhibits the lowest Oracle score (40.08%), indicating fundamental challenges in clustering feature sets meaningfully. Reinlection Accuracy + Hierarchical K-Means shows a particularly large Oracle-Algorithmic gap (65.64% \rightarrow 42.43%), highlighting that while clustering is effective, principal parts selection still requires refinement.

These findings emphasize the importance of optimizing both clustering effectiveness and principal parts selection to bridge the gap between Oracle and Algorithmic performance.

6.3 Interplay Between Transposition and Oracle Performance

Table 5 presents the impact of transposition on Reinlection Accuracy models under both Oracle and Algorithmic evaluations.

The results indicate that while transposition improves Oracle performance for Hierarchical K-Means (65.64% \rightarrow 67.70%), it has a negligible effect on Algorithmic scores, indicating that while transposition enhances clustering under ideal conditions, it does not meaningfully improve principal

parts selection. Additionally, Affinity Propagation exhibits sensitivity to data orientation, showing a slight decline in Oracle performance (58.78% → 58.51%), suggesting that its clustering mechanism relies on specific directional patterns that transposition may disrupt. Conversely, Hierarchical K-Means benefits from transposed data, likely due to its iterative refinement of clusters. However, since Algorithmic scores remain largely unchanged across models, these findings reinforce that refining selection heuristics, rather than adjusting data orientation, is the key to improving model performance.

7 Related Work

Early computational approaches to paradigm completion predominantly relied on the lemma as the central reference form, treating it as the sole input for generating full inflectional paradigms (Durrett and DeNero, 2013; Hulden, 2014; Nicolai et al., 2015; Ahlberg et al., 2015; Faruqui et al., 2016). However, Cotterell et al. (2017) highlighted the limitations of this approach, noting that forcing transformations to pass exclusively through the lemma can introduce unnecessary complexity. Instead, more flexible models leveraging multiple inflected forms have been proposed, allowing transformations to occur directly or via intermediary forms, rather than constraining them to a single privileged form. This shift aligns with the concept of principal parts, which constitute the minimal set of paradigm cells required to deduce all others (Finkel and Stump, 2007; Stump and Finkel, 2013).

Cotterell et al. (2017) introduced a directed graphical model that probabilistically generates missing inflected forms by modeling dependencies within paradigms. This approach enables the prediction of a form from multiple inflected forms rather than exclusively from the lemma. Around the same time, Kann et al. (2017) introduced multi-source reinflection, demonstrating that using multiple inflected forms as input improves accuracy. Their work explicitly references principal parts as a linguistic motivation, reinforcing the idea that certain forms within a paradigm hold stronger predictive capacity. Additionally, Cotterell et al. (2019) examined the structural complexity of inflectional paradigms, proposing a neural method for ordering paradigm slots based on their predictability—an indirect computational realization of the principal parts concept.

Liu and Hulden (2020) extended these ideas by reformulating morphological inflection as a Paradigm Cell Filling Problem (PCFP), where missing forms are inferred from a partially observed set of paradigm cells. While their work does not explicitly model principal parts, it aligns with their predictive role in improving inflectional accuracy, particularly in low-resource settings.

Despite these advancements, no prior work has proposed a systematic, data-driven approach to principal parts detection. Existing studies have either assumed pre-defined principal parts or incorporated them indirectly within broader inflectional tasks. In contrast, we introduce PRINCIPAL PARTS DETECTION as a formal computational task, developing a multilingual benchmark and a principled methodology for automatic PRINCIPAL PARTS DETECTION. By integrating linguistic insights with computational modeling, we establish a structured framework for principal parts detection.

8 Conclusions

This work introduces PRINCIPAL PARTS DETECTION as a computational task, formalizing the detection of principal parts within inflectional paradigms. We construct a multilingual dataset with ten typologically diverse languages, and develop a structured framework to automatically detect principal parts in their verbal diagrams.

our empirical evaluation shows that characterizing inter-cell relationships, clustering feature sets, and selecting representatives, offers a viable strategy for identifying principal parts. Our best-performing approach — Edit Script similarity with Hierarchical K-Means — achieves an F1 score of 55.05%, significantly surpassing the Random Baseline of 21.20%. However, results across models indicate that while clustering is effective in grouping related feature sets, principal parts selection remains a key bottleneck.

Beyond theoretical interest, solving PRINCIPAL PARTS DETECTION has practical implications for computational morphology. By identifying compact, generative subsets of paradigm forms, principal parts can be leveraged to optimize morphological inflection models, reduce annotation costs, and improve low-resource language modeling. The structured approach presented here lays the foundation for future advancements, underscoring the relevance of linguistic principles in shaping more efficient NLP methodologies.

722 Limitations

723 Despite the progress demonstrated in this study,
724 several open challenges remain. Irregular
725 paradigms, as seen in Latin, continue to pose dif-
726 ficulties, highlighting the need for methods that
727 can better capture morphological unpredictability.
728 Additionally, our reliance on UniMorph, while of-
729 fering broad linguistic coverage, exposes inconsis-
730 tencies that impact model generalization. More
731 curated linguistic resources could improve dataset
732 reliability and refine the evaluation of principal
733 parts across languages.

734 Also, one could explore alternative clustering
735 strategies that are better suited to morphological
736 structures, such as graph-based methods or neural
737 clustering approaches. Transformer-based models
738 hold potential for capturing deeper morphological
739 dependencies, offering an avenue for enhancing
740 both clustering accuracy and principal parts selec-
741 tion. These challenges are beyond the scope of this
742 paper and we reserve it to future work.

743 Our dataset currently includes only 10 languages.
744 Expanding the dataset to include more morpholog-
745 ically rich and underrepresented languages, such
746 as polysynthetic languages, would better capture
747 typological diversity and will potentially further
748 validate the robustness of PRINCIPAL PARTS DE-
749TECTION methods.

750 References

751 Malin Ahlberg, Markus Forsberg, and Mans Hulden.
752 2015. [Paradigm classification in supervised learning of morphology](#). In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1024–1029, Denver, Colorado. Association for Computational Linguistics.

759 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-
760 gio. 2014. Neural machine translation by jointly
761 learning to align and translate. *arXiv preprint arXiv:1409.0473*.

763 Khuyagbaatar Batsuren, Omer Goldman, Salam Khal-
764 ifa, Nizar Habash, Witold Kieraś, Gábor Bella,
765 Brian Leonard, Garrett Nicolai, Kyle Gorman, Yustinus Ghanggo Ate, Maria Ryskina, Sabrina Mielke,
766 Elena Budianskaya, Charbel El-Khaissi, Tiago Pimentel, Michael Gasser, William Abbott Lane,
767 Mohit Raj, Matt Coler, Jaime Rafael Montoya Samame, Delio Siticonatzi Camaiteri, Esaú Zumaeta Rojas, Didier López Francis, Arturo Oncevay, Juan López Bautista, Gema Celeste Silva Villegas, Lucas Torroba Hennigen, Adam Ek, David

Guriel, Peter Dirix, Jean-Philippe Bernardy, Andrey Scherbakov, Aziyana Bayyr-ool, Antonios Anastasopoulos, Roberto Zariquiey, Karina Sheifer, Sofya Ganieva, Hilaria Cruz, Ritván Karahóga, Stella Markantonatou, George Pavlidis, Matvey Plugaryov, Elena Klyachko, Ali Salehi, Candy Angulo, Jatayu Baxi, Andrew Krizhanovsky, Natalia Krizhanovskaya, Elizabeth Salesky, Clara Vania, Sardana Ivanova, Jennifer White, Rowan Hall Maudslay, Josef Valvoda, Ran Zmigrod, Paula Czarnowska, Irene Nikkarinen, Aelita Salchak, Brijesh Bhatt, Christopher Straughn, Zoey Liu, Jonathan North Washington, Yuval Pinter, Duygu Ataman, Marcin Wolinski, Totok Suhardijanto, Anna Yablonskaya, Niklas Stoehr, Hossep Dolatian, Zahroh Nuriah, Shyam Ratan, Francis M. Tyers, Edoardo M. Ponti, Grant Aiton, Aryaman Arora, Richard J. Hatcher, Ritesh Kumar, Jeremiah Young, Daria Rodionova, Anastasia Yemelina, Taras Andrushko, Igor Marchenko, Polina Mashkovtseva, Alexandra Serova, Emily Prud’hommeaux, Maria Nepomniashchaya, Fausto Giunchiglia, Eleanor Chodroff, Mans Hulden, Miikka Silfverberg, Arya D. McCarthy, David Yarowsky, Ryan Cotterell, Reut Tsarfaty, and Ekaterina Vylomova. 2022. [UniMorph 4.0: Universal Morphology](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 840–855, Marseille, France. European Language Resources Association.

Tadeusz Caliński and Jerzy Harabasz. 1974. A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods*, 3(1):1–27.

Ryan Cotterell, Christo Kirov, Mans Hulden, and Jason Eisner. 2019. On the complexity and typology of inflectional morphological systems. *Transactions of the Association for Computational Linguistics*, 7:327–342.

Ryan Cotterell, Christo Kirov, John Sylak-Glassman, David Yarowsky, Jason Eisner, and Mans Hulden. 2016. The sigmorphon 2016 shared task—morphological reinflection. In *Proceedings of the 14th SIGMORPHON workshop on computational research in phonetics, phonology, and morphology*, pages 10–22.

Ryan Cotterell, John Sylak-Glassman, and Christo Kirov. 2017. [Neural graphical models over strings for principal parts morphological paradigm completion](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 759–765, Valencia, Spain. Association for Computational Linguistics.

Greg Durrett and John DeNero. 2013. [Supervised learning of complete morphological paradigms](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1185–1195, Atlanta, Georgia. Association for Computational Linguistics.

774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832

833	Manaal Faruqui, Yulia Tsvetkov, Graham Neubig, and	Robert A Wagner and Michael J Fischer. 1974. The	888
834	Chris Dyer. 2016. Morphological inflection genera-	string-to-string correction problem. <i>Journal of the</i>	889
835	tion using character sequence to sequence learning.	<i>ACM (JACM)</i> , 21(1):168–173.	890
836	<i>arXiv preprint arXiv:1512.06110</i> .		
837	Raphael Finkel and Gregory Stump. 2007. Principal	Appendix	891
838	parts and morphological typology. <i>Morphology</i> ,	A Technical Overview of the PRINCIPAL	892
839	17:39–75.	PARTS DETECTION Dataset	893
840	Brendan J. Frey and Dmitri Dueck. 2007. Clustering	This section provides the technical details of the	894
841	by passing messages between data points. <i>Science</i> ,	PRINCIPAL PARTS DETECTION dataset, including	895
842	315(5814):972–976.	the number of samples per feature set in each lan-	896
843	Omer Goldman, Khuyagbaatar Batsuren, Salam Khal-	guage’s verb paradigm and the total number of gold	897
844	ifa, Aryaman Arora, Garrett Nicolai, Reut Tsarfaty,	principal parts for each language. In some cases,	898
845	and Ekaterina Vylomova. 2023. SIGMORPHON–	specific feature sets were removed for various rea-	899
846	UniMorph 2023 shared task 0: Typologically di-	sons, which are explained in subsection A.2 .	900
847	verse morphological inflection. In <i>Proceedings of the</i>	Additionally, we list the gold principal parts for	901
848	<i>20th SIGMORPHON workshop on Computational</i>	each language, formatted as <code>feature_set</code> (e.g.,	902
849	<i>Research in Phonetics, Phonology, and Morphology</i> ,	form). When two feature sets share the same form,	903
850	pages 117–125, Toronto, Canada. Association for	the gold principal parts are listed in square brackets	904
851	Computational Linguistics.	[]. The first feature set corresponds to the princi-	905
852	Omer Goldman, David Guriel, and Reut Tsarfaty. 2021.	pal part identified in linguistic literature, while the	906
853	(un) solving morphological inflection: Lemma over-	second represents a feature set that consistently	907
854	lap artificially inflates models’ performance. <i>arXiv</i>	shares the same form across all samples in the	908
855	<i>preprint arXiv:2108.05682</i> .	dataset. In such cases, the second feature set is	909
856	Mans Hulden. 2014. Generalizing inflection tables into	included as a possible principal part, as the algo-	910
857	paradigms with finite state operations. In <i>Proceed-</i>	rithm’s choice between them does not affect the	911
858	<i>ings of the 2014 Joint Meeting of SIGMORPHON</i>	analysis. To avoid redundancy, no principal part is	912
859	<i>and SIGFSM</i> , pages 29–36, Baltimore, Maryland.	counted more than once in these scenarios.	913
860	Association for Computational Linguistics.	A.1 Dataset Summary and Illustrative	914
861	Katharina Kann, Ryan Cotterell, and Hinrich Schütze.	Lexeme Examples	915
862	2017. Neural multi-source morphological reinflec-	For each language, we provide an example lex-	916
863	tion. In <i>Proceedings of the 15th Conference of the</i>	eme to illustrate the principal parts, formatted as	917
864	<i>European Chapter of the Association for Computa-</i>	<code>feature_set</code> (e.g., form). These examples are	918
865	<i>tional Linguistics: Volume 1, Long Papers</i> , pages	illustrative and may not share the same meanings	919
866	514–524, Valencia, Spain. Association for Computa-	across languages.	920
867	tional Linguistics.	A.2 Explanatory Notes	921
868	Vladimir I. Levenshtein. 1966. Binary codes capable of	The following explanatory notes clarify decisions	922
869	correcting deletions, insertions, and reversals. <i>Soviet</i>	made during dataset preparation and supplement	923
870	<i>physics doklady</i> , 10(8):707–710.	the information presented in Table 6 :	924
871	Ling Liu and Mans Hulden. 2020. Leveraging principal	• Spanish: PRO feature sets, representing verbs	925
872	parts for morphological inflection. In <i>Proceedings</i>	with object clitic pronouns, were removed.	926
873	<i>of the 17th SIGMORPHON Workshop on Computa-</i>	• Swedish: The V-IMP-PASS feature set was	927
874	<i>tional Research in Phonetics, Phonology, and Mor-</i>	excluded due to insufficient samples (only	928
875	<i>phology</i> , pages 153–161.	three).	929
876	Eugene W Myers. 1986. An o (nd) difference algorithm	• Latin:	930
877	and its variations. <i>Algorithmica</i> , 1(1):251–266.	– Passive feature sets were excluded.	931
878	Garrett Nicolai, Colin Cherry, and Grzegorz Kondrak.	– Feature sets starting with V.PTCP (in-	932
879	2015. Inflection generation as discriminative string	stead of V-V.PTCP) were removed.	933
880	transduction. In <i>Proceedings of the 2015 Conference</i>		
881	<i>of the North American Chapter of the Association</i>		
882	<i>for Computational Linguistics: Human Language</i>		
883	<i>Technologies</i> , pages 922–931, Denver, Colorado. As-		
884	sociation for Computational Linguistics.		
885	Gregory Stump and Raphael A. Finkel. 2013. <i>Principal</i>		
886	<i>parts</i> , Cambridge Studies in Linguistics, page 9–39.		
887	Cambridge University Press.		

Language	Features	Samples per Feature Set	# of Gold Principal Parts	Gold Principal Parts
English	5	23,896–31,848	3	V-NFIN-IMP+SBJV (e.g., eat), V-PST (e.g., ate), V-V.PTCP-PST (e.g., eaten)
French	49	7,483–7,535	7	V-NFIN (e.g., mangier), V-IND-PRS-1-PL (e.g., manjons), V.PTCP-PST (e.g., mangié), V-IND-FUT-1-SG (e.g., mangerai), V-IND-PRS-1-SG (e.g., manju), V-IND-PRS-3-PL (e.g., manjäent), V-IND-PST-1-SG-PFV (e.g., manjai)
German	30	2,307–6,661	5	V-NFIN (e.g., essen), V.PTCP-PST (e.g., gegessen), [V-IND-SG-3-PST, V-IND-SG-1-PST (e.g., aß)], V-IND-SG-3-PRS (e.g., isst), [V-SBJV-SG-3-PST, V-SBJV-SG-1-PST (e.g., äße)]
Spanish	79	6,676–6,695	2	V-NFIN (e.g., comer), V-IND-PRS-1-SG (e.g., como)
Danish	8	162	5	V-ACT-NFIN (e.g., danse), V-ACT-IND-PRS (e.g., danser), V-ACT-IND-PST (e.g., dansede), V-ACT-IMP (e.g., dans), V.PTCP-PASS-PST (e.g., danset)
Swedish	19	2,114–2,536	5	[V-NFIN-ACT, V-IND-PL-ACT-PRS (e.g., äta)], V-IND-SG-ACT-PRS (e.g., äter), V-IND-SG-ACT-PST (e.g., ät), V-V.CVB-ACT (e.g., ätit), V-IMP-ACT (e.g., ät)
Finnish	161	7,221–7,226	4	V-NFIN-ACT+PASS (e.g., syödä), V-ACT-PRS-POS-IND-1-SG (e.g., syön), V-ACT-PST-POS-IND-3-SG (e.g., söi), V.PTCP-ACT-PST (e.g., syönyt)
Turkish	703	588	2	V-NFIN (e.g., içmek), V-IND-PRS-HAB-3-SG-POS-DECL (e.g., içer)
Latin	48	450–947	3	V-IND-ACT-PRS-1-SG (e.g., -plēō), V-NFIN-ACT-PRS (e.g., -plēre), V-V.MSDR-ACC-LGSPEC1 (e.g., -plētum)

Table 6: Summary of the PRINCIPAL PARTS DETECTION dataset by language, including gold principal parts and illustrative lexeme examples.

- Feature sets with 30 or fewer samples were excluded.

B Characterization Tables for Selected Languages

- The first-person-singular-perfect-active-indicative feature set was excluded from the gold principal parts list due to insufficient data (only two samples).

To illustrate the structure of the characterization methods, we present detailed characterization tables for three representative languages from our dataset. These tables demonstrate how different feature sets relate within their verb paradigms, showcasing the variation across Edit Distance, Edit Script, and Reinflection Accuracy characterization methods.

950 Each language is represented by three tables, cor-
951 responding to the distinct characterization methods,
952 with principal parts highlighted in yellow for clarity.
953 Additionally, cases where two feature sets consis-
954 tently share the same form and are interchangeable
955 as principal parts are marked with a distinct color.
956 Since these feature sets carry identical information,
957 the model’s selection between them does not im-
958 pact the results.

959 **Interpretation of Tables.** The provided tables
960 exemplify the structure of the characterization
961 methods rather than an exhaustive display of all
962 ten languages in our study. While specific lexeme
963 examples are shown in the rows and columns, the
964 quantified relationships they capture apply to the
965 entire verb paradigm of each language. These ex-
966 amples serve to illustrate the broader implicative
967 patterns identified during the characterization pro-
968 cess.

969 **B.1 Characterization Tables for English**

970 **B.2 Characterization Tables for German**

971 **B.3 Characterization Tables for Swedish**

	Features	V-NFIN-IMP+SBJV - eat	V-PRS-3-SG - eats	V-PST - ate	V-V.PTCP-PRS - eating	V-V.PTCP-PST - eaten
1	V-NFIN-IMP+SBJV - eat	0	1.157683294	1.532683294	3.088508537	1.534943087
2	V-PRS-3-SG - eats	1.157683294	0	1.421493137	3.087504185	1.410905591
3	V-PST - ate	1.532683294	1.421493137	0	3.066078005	0.048627385
4	V-V.PTCP-PRS - eating	3.088508537	3.087504185	3.066078005	0	3.034273519
5	V-V.PTCP-PST - eaten	1.534943087	1.410905591	0.048627385	3.034273519	0

Figure 1: Average edit distances for the English verb paradigm. Values range from 0 to 4.5. Darker red shades indicate closer relationships between feature sets, while darker turquoise shades represent greater differences.

	Features	V-NFIN-IMP+SBJV - eat	V-PRS-3-SG - eats	V-PST - ate	V-V.PTCP-PRS - eating	V-V.PTCP-PST - eaten
1	V-NFIN-IMP+SBJV - eat	1	27	117	51	124
2	V-PRS-3-SG - eats	29	1	110	48	117
3	V-PST - ate	124	110	1	116	43
4	V-V.PTCP-PRS - eating	55	59	119	1	121
5	V-V.PTCP-PST - eaten	128	118	45	119	1

Figure 2: Edit Script counts for the English verb paradigm. Values range from 1 to 128. Darker purple shades indicate fewer unique character sets (closer relationships), while darker air-force-blue shades reflect greater variation.

	Features	V-NFIN-IMP+SBJV - eat	V-PRS-3-SG - eats	V-PST - ate	V-V.PTCP-PRS - eating	V-V.PTCP-PST - eaten
1	V-NFIN-IMP+SBJV - eat	0.95	0.96	0.92	0.94	0.92
2	V-PRS-3-SG - eats	0.95	0.96	0.92	0.94	0.91
3	V-PST - ate	0.9	0.91	0.96	0.94	0.95
4	V-V.PTCP-PRS - eating	0.91	0.92	0.92	0.95	0.92
5	V-V.PTCP-PST - eaten	0.91	0.91	0.96	0.95	0.96

Figure 3: Reinfection Accuracy scores for the English verb paradigm. Values range from 0.9 to 0.96. Darker teal shades indicate higher accuracy, while darker pink shades reflect lower performance.

Feature	V-NFIN-IMP+SBJV - eat	V-PRS-3-SG - eats	V-PST - ate	V-V.PTCP-PRS - eating	V-V.PTCP-PST - eaten
V-NFIN-IMP+SBJV - eat	0.95	0.96	0.92	0.94	0.92
V-PRS-3-SG - eats	0.95	0.96	0.92	0.94	0.91
V-PST - ate	0.9	0.91	0.96	0.94	0.95
V-V.PTCP-PRS - eating	0.91	0.92	0.92	0.95	0.92
V-V.PTCP-PST - eaten	0.91	0.91	0.96	0.95	0.96

Figure 4: Average edit distances for the German verb paradigm. Values range from 0 to 11.19. Darker red shades indicate closer relationships between feature sets, while darker ball-blue shades represent greater distances.

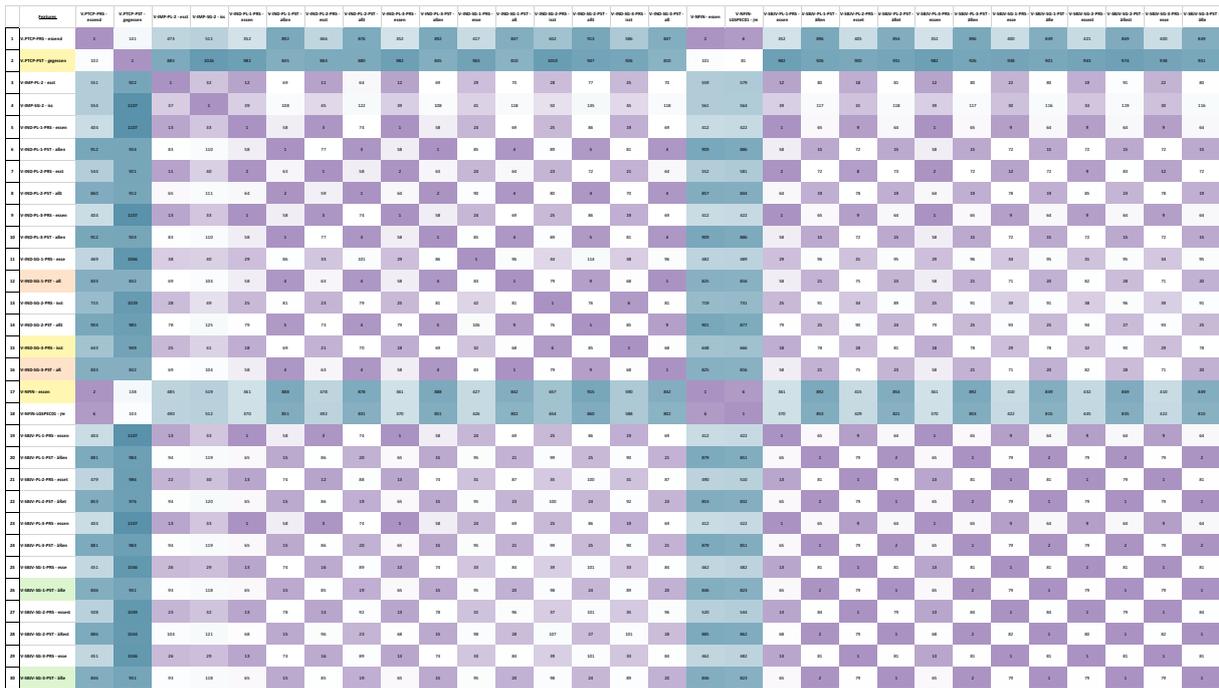


Figure 5: Edit Script scores for the German Verb Paradigm. Values range from 1 to 1,107. Darker purple shades indicate fewer unique character sets (closer relationships), while darker air-force-blue shades reflect greater variation.

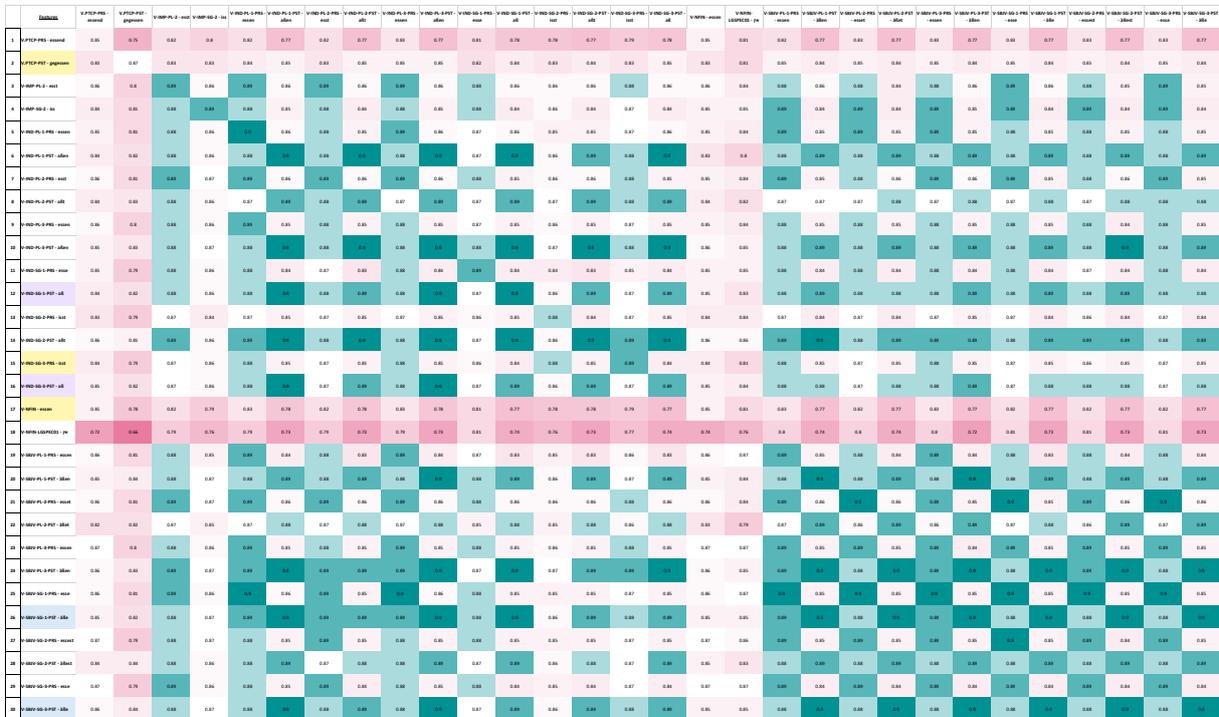


Figure 6: Reinfection Accuracy scores for the German verb paradigm. Values range from 0.66 to 0.9. Darker teal shades indicate higher accuracy, while darker pink shades reflect lower performance.

	Feature	V-MP-ACT - It	V-IND-PL-ACT-PRES - Ita	V-IND-PL-ACT-PST - Ita	V-IND-PL-PASS-PRES - Ita	V-IND-PL-PASS-PST - Ita	V-IND-SG-ACT-PRES - It	V-IND-SG-ACT-PST - It	V-IND-SG-PASS-PRES - It/Ita	V-IND-SG-PASS-PST - Ita	V-MFIN-ACT - Ita	V-MFIN-PASS - Ita	V-SBV-ACT-PRES - Ita	V-SBV-ACT-PST - Ita	V-SBV-PASS-PRES - Ita	V-SBV-PASS-PST - Ita	V-V-CVB-ACT - Ita	V-V-CVB-PASS - Ita	V-V-PTCP-PRES - Ita	V-V-PTCP-PST - Ita
1	V-MP-ACT - It	0	0.27297804	2.08505178	1.27522012	3.08120463	1.23558057	1.05147029	1.01657982	2.98352496	0.27240264	1.17125012	0.96449701	2.07495069	1.98261487	1.08965244	1.23679313	2.10280668	1.45414384	1.43098905
2	V-IND-PL-ACT-PRES - Ita	0.27297804	0	2.08783754	1	3.07578463	1.22949268	2.05717666	1.03944676	2.96878197	0	1	0.93228706	2.07728706	1.97386177	1.06444443	1.27011041	2.30537186	1.39639804	1.461606747
3	V-IND-PL-ACT-PST - Ita	2.08505178	2.08783754	0	3.16011369	1.00212081	2.17426011	0.17389899	2.06142956	1.03789671	2.08001017	2.116011369	1.02152968	0.12815474	2.82210564	1.11747887	2.02865615	2.30037897	2.00050495	1.46501542
4	V-IND-PL-PASS-PRES - Ita	1.27522012	1	2.216011369	0	2.08781826	1.26382283	2.18715111	0.27483437	2.05806575	1	0	1.99261014	2.16011369	0.93931827	2.07615894	1.43028612	1.29202703	3.1974874	1.61431316
5	V-IND-PL-PASS-PST - Ita	3.08120463	3.07578463	1.00521081	2.08781826	0	2.93229589	1.15046798	2.0606187	0.14522327	0.97860339	2.08781826	2.90529817	1.13610449	1.90881171	0.11236781	2.98878144	2.07852415	2.84297208	2.32148878
6	V-IND-SG-ACT-PRES - Ita	1.23558057	1.23495268	2.17426011	1.26382283	1.93229589	0	2.14037849	1.30212643	2.9052108	1.23887115	1.26382283	1.08611887	2.08670789	1.71361757	2.84320174	1.41284313	2.10133955	3.08270219	1.54005263
7	V-IND-SG-ACT-PST - It	1.05147029	2.07717666	0.17389899	2.18715111	1.50505798	2.14037849	0	2.0568782	1.00521081	2.05721784	2.18715111	1.89072871	0.17389899	2.83254382	1.15016798	2.10755511	2.21318914	2.01108910	1.485130917
8	V-IND-SG-PASS-PRES - It/Ita	1.03944676	1.03843676	2.09616296	0.27483437	2.09061887	2.18022643	2.0568782	0	1.95817753	1.020397819	0.27483437	1.7074775	2.09542865	0.93223273	2.052986132	1.25817483	1.30955165	1.45480297	1.426512968
9	V-IND-SG-PASS-PST - Ita	2.98352496	2.96878197	1.03789671	2.05806575	0.14522327	2.90052108	1.00521081	1.96877751	0	2.96863864	2.05806575	2.78725487	1.03894813	1.87795645	0.14522327	2.98024216	2.0884479	2.83388547	2.32177185
10	V-MFIN-ACT - Ita	0.27240264	0	2.08803157	1	3.07585339	1.23887115	2.05721784	1.020397819	2.96863864	0	1	0.93224948	2.07734806	1.97525747	1.06446538	1.26953413	2.05784732	1.39639206	1.46112887
11	V-MFIN-PASS - Ita	1.715225012	1	2.216011369	0	2.08781826	1.26382283	2.18715111	0.27483437	2.05806575	1	0	1.99261014	2.16011369	0.93931827	2.07615894	1.43028612	1.29202703	3.1974874	1.61431316
12	V-SBV-ACT-PRES - Ita	0.96449701	0.93228706	1.92152968	1.95026014	2.90529817	1.03911887	1.89072871	1.70740775	2.78725487	0.93224948	1.93920104	0	1.82018214	1	2.81477975	1.96960145	2.01942207	3.22058801	2.05814686
13	V-SBV-ACT-PST - Ita	2.07495069	2.07717666	0.17389899	2.18715111	1.50505798	2.08670789	0.17389899	2.00424865	1.00504913	2.07734806	2.16011369	1.82018214	0	2.78027667	1.00426382	2.02865615	2.00037897	1.89277228	1.351512302
14	V-SBV-PASS-PRES - Ita	1.96813487	1.97386177	2.30037897	0.95818827	1.90681171	1.72382757	2.83134382	0.93223273	1.87795641	1.95132947	0.93931827	1	2.78017667	0	1.814560136	1.23976905	1.94613772	4.13089349	2.23150845
15	V-SBV-PASS-PST - Ita	1.08965244	1.06444443	1.11747887	2.07815894	0.11236781	2.84320274	1.15050578	2.05280112	0.14522327	0.06446538	2.07615894	2.81477975	1.00430182	1.814849538	0	2.98878144	2.07852415	2.84297208	2.32070001
16	V-V-CVB-ACT - Ita	1.23621913	1.27011041	2.09265615	1.43828612	2.98878144	1.41841633	2.10720551	1.26817463	2.980504216	1.26953413	1.45818612	1.98678144	2.02865615	2.129766305	2.98878144	0	1.01421127	1.58217818	1.349312028
17	V-V-CVB-PASS - Ita	2.12800866	2.30537186	2.20037897	1.212601703	2.07852415	2.25153955	2.21338914	1.30955165	2.0848479	2.208784732	1.29202703	2.90042207	2.20037897	1.94512772	2.07852415	1.003421127	0	3.57095103	2.38818448
18	V-V-PTCP-PRES - Ita	1.45414384	1.39639804	2.00050495	3.1974874	2.84297208	3.48172673	2.01108910	3.45480297	2.833885471	3.19619206	3.1974874	3.22058801	1.89277228	4.13089349	2.814591732	3.58217818	3.57095103	0	2.46758489
19	V-V-PTCP-PST - Ita	1.43098905	1.461606747	1.40560142	1.61431316	2.32148878	1.54005263	1.485130917	1.48512968	2.32177185	1.46112887	1.61431316	2.02814686	1.351512302	2.21208105	2.23670001	1.949312028	2.38818448	2.46758489	0

Figure 7: Average edit distances for the Swedish verb paradigm. Values range from 0 to 4.153. Darker red shades indicate closer relationships between feature sets, while darker blue shades represent greater differences.

	Feature	V-MP-ACT - It	V-IND-PL-ACT-PRES - Ita	V-IND-PL-ACT-PST - Ita	V-IND-PL-PASS-PRES - Ita	V-IND-PL-PASS-PST - Ita	V-IND-SG-ACT-PRES - It	V-IND-SG-ACT-PST - It	V-IND-SG-PASS-PRES - It/Ita	V-IND-SG-PASS-PST - Ita	V-MFIN-ACT - Ita	V-MFIN-PASS - Ita	V-SBV-ACT-PRES - Ita	V-SBV-ACT-PST - Ita	V-SBV-PASS-PRES - Ita	V-SBV-PASS-PST - Ita	V-V-CVB-ACT - Ita	V-V-CVB-PASS - Ita	V-V-PTCP-PRES - Ita	V-V-PTCP-PST - Ita
1	V-MP-ACT - It	1	6	52	7	54	6	47	5	48	6	7	6	51	6	53	37	38	66	80
2	V-IND-PL-ACT-PRES - Ita	6	1	55	1	51	5	61	7	55	1	1	7	54	6	50	36	34	48	88
3	V-IND-PL-ACT-PST - Ita	54	54	1	49	4	56	11	57	13	54	49	57	2	52	5	33	34	116	70
4	V-IND-PL-PASS-PRES - Ita	7	1	51	1	51	6	55	7	55	1	1	7	50	6	50	35	34	33	79
5	V-IND-PL-PASS-PST - Ita	55	50	4	49	1	54	13	54	10	50	49	52	5	52	2	34	31	200	67
6	V-IND-SG-ACT-PRES - Ita	6	6	63	6	59	1	65	8	61	6	6	7	53	6	57	39	38	65	83
7	V-IND-SG-ACT-PST - It	47	55	11	50	13	57	1	50	4	55	50	58	11	53	13	40	41	124	77
8	V-IND-SG-PASS-PRES - It/Ita	5	8	55	7	53	7	49	1	47	8	7	6	54	5	52	36	36	51	76
9	V-IND-SG-PASS-PST - Ita	48	51	13	50	10	56	4	47	1	51	50	53	13	53	10	38	38	103	72
10	V-MFIN-ACT - Ita	6	1	55	1	51	5	61	7	55	1	1	7	54	6	50	36	34	48	88
11	V-MFIN-PASS - Ita	7	1	51	1	51	6	55	7	55	1	1	7	50	6	50	35	34	31	79
12	V-SBV-ACT-PRES - Ita	6	7	64	6	57	6	66	5	58	7	6	5	54	1	50	35	32	62	88
13	V-SBV-ACT-PST - Ita	51	53	2	48	5	55	11	56	13	53	48	56	1	51	4	31	34	124	71
14	V-SBV-PASS-PRES - Ita	7	7	56	6	57	6	57	5	58	7	6	1	49	1	50	34	33	47	79
15	V-SBV-PASS-PST - Ita	54	49	5	48	2	53	13	53	10	49	48	51	4	51	1	34	31	104	68
16	V-V-CVB-ACT - Ita	38	40	33	36	34	40	39	38	39	40	36	40	33	35	34	1	1	100	60
17	V-V-CVB-PASS - Ita	38	37	36	36	32	40	41	36	37	37	36	39	35	36	32	3	1	88	54
18	V-V-PTCP-PRES - Ita	66	48	137	33	115	68	163	57	119	48	33	64	128	48	108	106	91	1	39
19	V-V-PTCP-PST - Ita	81	90	80	79	74	91	85	78	79	80	79	87	73	76	68	56	40	34	1

Figure 8: Edit Script scores for the Swedish verb paradigm. Values range from 1 to 142. Darker purple shades indicate fewer unique character sets (closer relationships), while darker air-force-blue shades reflect greater variation.

