

A Computational Framework to Identify Self-Aspects in Text

Anonymous ACL submission

Abstract

This Ph.D. proposal introduces a plan to develop a computational framework to identify Self-aspects in text. The Self is a multifaceted construct and it is reflected in language. While it is described across disciplines like cognitive science and phenomenology, it remains underexplored in natural language processing (NLP). Many of the aspects of the Self align with psychological and other well-researched phenomena (e.g., those related to mental health), highlighting the need for systematic NLP-based analysis. In line with this, we plan to introduce an ontology of Self-aspects and a gold-standard annotated dataset. Using this foundation, we will develop and evaluate conventional discriminative models, generative large language models, and embedding-based retrieval approaches against four main criteria: interpretability, ground-truth adherence, accuracy, and computational efficiency. Top-performing models will be applied in case studies in mental health and empirical phenomenology.

1 Introduction

The Self, superficially experienced as “the (perhaps sometimes elusive) feeling of being the particular person one is” (Siderits et al., 2013), is a complex phenomenon, amply discussed in philosophy and cognitive science (e.g., Zahavi, 2008). While there exist different views about the metaphysical nature of the Self (Siderits et al., 2013), in this work, we build on its phenomenological and behavioural manifestations. In everyday experience, the Self is characterised by multiple phenomenological and psychological aspects, including the experience of one’s own body (Bermúdez, 2018) and a sense of agency (Gallagher, 2000), among others (Caporusso, 2022).

These Self-aspects are conceptually and empirically related to other well-established constructs—such as personality traits or experiential modes. For example, their relevance to contexts

such as mental health research is supported in related work, which highlights the central role of Self-related processes in well-being and psychopathology, as well as in empirical phenomenology (i.e., the empirical investigation of experience, Aspers, 2009), where they are key to understanding altered states of consciousness (see Section 2).

Importantly, the specific ways in which Self-aspects are experienced by a person in a given moment are reflected in the language they use (e.g., see Pennebaker et al. (2003) and Section 2). The found correlations between textual features and Self-aspects can be further employed in downstream NLP tasks, for instance to detect psychological states (Caporusso et al., 2023; Du and Sun, 2022; Kolenik et al., 2024). However, the connections between textual features and many Self-aspects important for the identification of, e.g., mental health conditions and phenomenological states, are underexplored.

To address this shortcoming, we propose a computational framework capable of automatically detecting the presence and mode of Self-aspects in text. Existing tools such as LIWC (Linguistic Inquiry and Word Count; Boyd et al., 2022) and VADER (Valence Aware Dictionary and sEntiment Reasoner; Hutto and Gilbert, 2014) have shown that psychologically meaningful patterns can be computationally extracted from text using lexicons and interpretable features. Building on this tradition, our framework aims to go further: to detect nuanced, theoretically grounded aspects of Self-experience—such as agency, embodiment, or narrative coherence—through a combination of ontology design, annotated data, and a range of modelling approaches. The resulting method can be applied to tasks in domains such as mental health research and empirical phenomenology.

2 Related Work

2.1 Textual Features and Self-Aspects Correlations

This subsection surveys studies mapping text features to aspects of the Self.

Self Aspects Most research focuses on *I-talk*, i.e., the use of first-person pronouns as indicators of Self-focus (Pennebaker et al., 2003), which correlates with emotional pain, trauma, and depression (Tausczik and Pennebaker, 2010). Furthermore, pronoun usage hints at specific understandings of the Self vs others distinction (Na and Choi, 2009; Sharpless, 1985). The usage of active vs passive voice can shed light on the sense of agency of the author of a text (Simchon et al., 2023), while the Narrative Self (NS; i.e., “the narrative someone has of themselves, comprising their autobiographical memories and stories of who they are” Caporusso et al., 2024) is reflected in the structure and coherence of one’s autobiographical accounts (Habermas and Köber, 2015; Holm et al., 2016; Jaeger et al., 2014; Waters and Fivush, 2015). In this context, Author profiling (AP) refers to the task of inferring personal characteristics of an author based on their writing, which has applications in, e.g., sociolinguistics and mental health analytics (Eke et al., 2019; Ouni et al., 2023b).

The correlation of text features with other aspects of the Self, such as the Minimal Self (MS; “the fact that experiences are presented to us in a fundamentally personal and subjective way” Caporusso et al., 2024), are less explored (Uno and Imaizumi, 2025).

Caporusso et al. (2024) investigated the LIWC categories associated with different aspects of the Self: MS, NS, Self as Agent (AS; “the experience of being an agent, i.e., in control, active”), Bodily Self (BS; “the experience of owning, controlling, and/or identifying with someone’s own body (or parts of it)”), and Social Self (SS; “the self as it is shaped and/or perceived when in an interaction or relationship of sorts with other people or entities to whom we attribute qualities of an inner life”). Specifically, utilising a mixed approach to annotate the data, the authors classified text instances as presenting or not each of the mentioned self-aspects, and they analysed the obtained splits with LIWC.

Methods The methodological approaches utilised to detect correlations between textual

features and Self-aspects can be broadly grouped into three main types:

- Approaches based on stylistic features such as punctuation, syntactic patterns, part-of-speech (POS) tags, sentence length, character/word n-grams, and structural features (e.g., number of paragraphs or capitalised words)—see Ouni et al. (2021); Vijayan and Govilkar (2019).
- Content-based approaches, relying on subject matter and vocabulary; features include term frequency-inverse document frequency (TF-IDF), topic models, and domain-specific keywords—see Ch and Cheema (2018); Ouni et al. (2023b)
- Hybrid approaches, where both stylistic and content-based features are analysed—see Fatima et al. (2017); Ouni et al. (2021, 2023b)

The use of LIWC or other lexicon-based techniques is the most common approach to investigate correlations between Self-aspects and textual features (Boyd and Schwartz, 2021; Pennebaker et al., 2003). More recently, however, NLP research has increasingly adopted machine learning (ML) methods—such as topic modelling and supervised classification—to analyse language patterns in a data-driven way (Eichstaedt et al., 2018; Ouni et al., 2021). Many studies used classical supervised learning methods, like support vector machines (SVMs; Chinea-Rios et al., 2022; HaCohen-Kerner, 2022; Vijayan and Govilkar, 2019), random forests (RFs; Fatima et al., 2017; Ouni et al., 2021), decision trees (Vijayan and Govilkar, 2019), and Naïve Bayes (NB; Mechti et al., 2020). Feature extraction in AP is critical: common strategies include Bag-of-Words (BoW) and TF-IDF (Ouni et al., 2023b), character and word n-grams (HaCohen-Kerner, 2022), POS and syntactic feature vectors (Mechti et al., 2020; Vijayan and Govilkar, 2019), word embeddings (Chinea-Rios et al., 2022; Fatima et al., 2017), semantic graphs and emotion tags (Ouni et al., 2023b). Furthermore, many studies employ qualitative approaches (Habermas and Köber, 2015; Waters and Fivush, 2015). However, deep learning (DL) models are increasingly employed as well, due to their capacity to automatically learn hierarchical feature representations from raw text and their superior performance on large-scale NLP tasks (Ouni et al., 2023a). Transformer-based models such as BERT (Devlin et al., 2019)

and RoBERTa (Liu et al., 2019) were adapted to AP tasks by fine-tuning on labelled AP datasets (Chinea-Rios et al., 2022). In recent work, LLMs have been explored for AP (see Huang et al., 2025). Huang et al. (2024) show that GPT-4 outperforms BERT-based models in zero-shot authorship attribution and verification, especially when guided by linguistic cues.

The type of text analysed varies widely, ranging from autobiographical essays (Adler, 2012; McAdams, 2001), stream-of-consciousness essays or narrative prompts (Pennebaker and Beall, 1986; Rude et al., 2004), transcripts of spoken conversations or interviews (Adler et al., 2008; Bamberg, 2008; Lysaker and Lysaker, 2002), diary entries and letters (Baumeister et al., 1994; Pennebaker and Francis, 1996), social media posts (Guntuku et al., 2019; Schwartz et al., 2013), to even published autobiographies or literature (Bruner, 2003; Freeman, 2009).

2.2 Downstream Applications

The correlations discussed in the previous subsection are often employed in downstream applications. For instance, Kolenik et al. (2024) utilised predefined sets of words and linguistic patterns that have been associated with specific psychological states, traits, or cognitive processes to train ML models that detect stress, anxiety, and depression. Similarly, Du and Sun (2022) leveraged linguistic features known to correlate with psychological states, like absolutist words and personal pronouns, to detect depression, anxiety, and suicidal ideation. In the context of the LT-EDI@RANLP 2023 shared task, first-person singular pronouns and time-related terms, recognised as indicative of depressive states (Ratcliffe, 2014), were employed to identify signs of depression in social media posts (Caporusso et al., 2023). Eichstaedt et al. (2018) utilised topic models to identify clusters of words that often appear together in Self-narratives, and supervised ML to predict an upcoming depression diagnosis from social media posts.

Outside of the context of NLP studies, works investigating, e.g., mental health issues or phenomenological states vastly address Self-aspects to identify the phenomenon of interest. For instance, an impacted sense of agency is registered in individuals with anxiety and depression, who experience a deficiency in estimating their control over positive outcomes (Mehta et al., 2023), while disturbances in interoception and Self-awareness were found

to be correlated with anxiety and schizophrenia, among the others (Yang et al., 2024). Often, different Self-aspects correlate with disorders in a synergistic way, or there is an atypical disintegration of Self-aspects. For instance, Alzheimer’s Disease (AD) and other conditions involving cognitive decline are associated with impaired Self-continuity, sense of personal history and future goals, capabilities of Self-reflection, and personal meaning (El Haj et al., 2015), resulting in a distorted narrative Self-identity. Along, and sometimes in support of, research in mental well-being, Self-aspects are relevant in the context of empirical phenomenology, among others. For example, a multitude of Self-aspects is examined in the investigation of experiences of dissolution (i.e., "experiential episodes during which the perceived boundaries between self and world (i.e., nonself) become fainter or less clear"; Caporusso, 2022; Nave et al., 2021), and bodily experience is investigated in the context of depersonalisation and derealisation disorders (Tanaka, 2018). In line with this, scales and symptom checklists have been developed to assess the presence and intensity of psychological or phenomenological states (Heering et al., 2016; Michal et al., 2014; Nour et al., 2016; Parnas et al., 2005; Sierra and Berrios, 2000).

2.3 Identified Gaps and Research Motivation

Disciplines like cognitive science, phenomenology, and psychology identify many different aspects of the Self, but NLP studies a) have dealt with only a few superficial ones and b) have only employed basic techniques. Indeed, while NLP started to employ the correlation between Self-aspects and textual features in various downstream tasks, the Self-aspects employed in, e.g., mental health research and empirical phenomenology are more varied and nuanced. For this reason, we believe that it would be helpful to identify further and more detailed connections between Self-aspects and textual features, and to develop a model to detect and analyse Self-aspects in text. This could be used by professionals of other disciplines, for instance to analyse patients’ reports and transcripts of phenomenological interviews (e.g., see micro-phenomenology, Petitmengin et al., 2019).

To this end, our proposed framework aligns in spirit with existing tools like LIWC and VADER. However, unlike these general-purpose approaches, our framework is specifically designed to capture a range of Self-aspects grounded in interdisciplinary

theory. Moreover, while LIWC captures psychological correlates at a coarse granularity (e.g., affect, pronouns), we aim to represent structured components of Self-experience.

3 Research Proposal

This Ph.D. proposal seeks to explore the ways of developing a computational model to automatically detect Self-aspects in language. We plan to test the proposed approaches on different case studies from the fields of mental health and empirical phenomenology. Our Research Objectives (ROs) are as following:

- **RO1)** Detail an ontology of the Self aspects that would be relevant and sensible for a computational model to detect in text.
- **RO2)** Construct heterogeneous datasets with annotations relative to the identified Self-aspects.
- **RO3)** Define the desiderata of the computational model to detect Self-aspects in text and identify the approaches which would best fulfill them.
- **RO4)** Determine the evaluation approach and the applications for our computational model to detect Self-aspects in text.

We plan to produce the following outcomes: self ontology detailing and labelling instructions; heterogeneous annotated dataset; set of models to identify Self-aspects in text.

4 Self Ontology (RO1)

We aim to develop a comprehensive ontology of Self-aspects which are a) relevant to possible applications and b) detectable in text data. Each Self-aspect (e.g., bodily Self) is characterised by different elements (e.g., body ownership, body awareness), each of which is specified in different modes (e.g., body ownership: weak). Some of the Self-aspects investigated are identified through previous studies which developed similar lists or ontologies (e.g., [Caporusso, 2022](#); [Nave et al., 2021](#)). The ontology, still a work-in-progress, is built collaboratively by adopting both bottom-up and a top-down approaches. That is to say, we utilise literature detailing the elements and modes of various Self-aspects (e.g., [Moore, 2016](#); [Serino et al., 2013](#))

along with studies from disciplines like psychology and neuroscience detailing the Self-aspects relevant to the construct of interest (e.g., [Petkova et al., 2011](#)). Furthermore, we will be meeting with experts from fields which could benefit from our final model (e.g., mental health professionals and empirical phenomenologists) to better identify the specific Self-aspects, elements, and modes which could be relevant for their work. While analysing literature and consulting with experts, we will be exploring textual data itself. For each Self-aspect, element, and mode, we will provide a definition, both a positive and a negative example from textual data, and notes to guide the identification and/or distinction. Constructing the Self ontology presents various challenges, most of all regarding how the different components relate with each other. For example, most of the aspects and elements, if not all, appear to not be mutually exclusive, and there are aspects (e.g., sense of agency) that could apply to other aspects (e.g., sense of agency over bodily Self).

5 Datasets (RO2)

The datasets (aiming for at least 10; see Section 8), which will be annotated with the labels developed (see Section 4), need to vary in type as it is desired for the model to be able to analyse Self-aspects across different kinds of data. We plan to utilise transcripts from phenomenological interviews, clinical tasks, and structured or unstructured interviews. These will include employ already existing datasets and construct new ones. Importantly, all data collection—whether previously conducted or ongoing—is carried out within the scope of pre-approved research projects. Part of the phenomenological interviews data has already been collected (six subjects), and clinical interviews are being conducted in the context of an existing larger project. We aim to utilise datasets from different languages, in order to create a multilingual model. The annotated datasets will serve as training and testing data, as well as ground truth. The length of the text chunk considered as a labelling instance is determined case by case, based on what is sufficient to meaningfully express the presence of a specific Self-aspect or mode. In general, this can range from a single sentence to a short paragraph, depending on the complexity of the expression.

5.1 Annotation

Multiple annotators (e.g., three, possibly the same researchers compiling the Self ontology and the annotation guidelines) will independently annotate the datasets or part of them. Inter-annotator agreement will be calculated to assess consistency and reliability of the annotations. The first author, who will take part in and lead the annotation, has experience in conducting qualitative analysis and annotation of textual data, including mostly phenomenological interviews, but also, e.g., social media posts, with a focus on the Self. In the first phase of the annotation process, the annotators will meet and discuss their decisions, so to come to a similar understanding of the guidelines. This can bring to further adjustments of the guidelines themselves. In the case that it proves too expensive to manually label the entire dataset, we will adopt large language models (LLMs) for automatic annotation of the remaining instances—following an approach similar to [Caporusso et al. \(2024\)](#). Specifically, LLMs fine-tuned for instruction following ([Brown et al., 2020](#)) will be evaluated against a manually annotated subset to ensure quality. Importantly, LLM-based annotations will be used to augment training data for conventional discriminative models, embedding-based retrieval approaches, and, in principle, for fine-tuning LLMs—provided such synthetic data is excluded from evaluation (see Section 7). LLMs themselves will be evaluated separately, using only the manually labelled portion of the data to avoid circularity. This ensures a clean separation between training supervision and model evaluation.

6 Desiderata (RO3a)

Here, we discuss our desiderata for the models.

Interpretability, which in the context of ML refers to the extent to which a human can understand the internal mechanism of a model leading from input to output ([Lipton, 2018](#); [Molnar, 2020](#)) is to be differentiated from explainability, which often involves post-hoc approximations of a model’s behaviour ([Molnar, 2020](#)). This distinction is particularly crucial for our task for three main reasons. First, the target applications of our framework include implementations in sensitive domains like healthcare. Indeed, in such cases, the use of interpretable ML models is preferable to post-hoc explanations for black-box models, as the latter may be incomplete or misleading and do not ensure

transparency, trust, and ethical decision-making ([Ahmad et al., 2018](#); [Amann et al., 2020](#); [Bohlen et al., 2024](#); [Chaddad et al., 2023](#); [Doshi-Velez and Kim, 2017](#); [Ennab and Mcheick, 2024](#); [Lipton, 2018](#); [Lu et al., 2023](#); [Rudin, 2019](#); [Tjoa and Guan, 2020](#)). Some examples of this are studies by [Gao et al. \(2023\)](#); [Wang et al. \(2023\)](#). Second, generic explainability approaches are often insufficient in NLP due to the inherent ambiguity, subjectivity, and domain sensitivity of language data, necessitating explanations that align with the linguistic and reasoning norms of specific application areas ([Mohammadi et al., 2025](#)). Some examples are studies by [Saha et al. \(2022, 2023\)](#); [Wang et al. \(2023\)](#). Third, interpretability is desirable because it enables traceability—the ability to identify the specific passage or linguistic marker that led to a given classification. This is particularly important in applications such as studies based on the analysis of empirical phenomenological interviews, where it is necessary to provide illustrative examples for each identified experiential category (e.g., a specific mode of a Self-aspect).

Ground-Truth Basis requires that model outputs be directly derived from verified, annotated data, rather than inferred through non-transparent or heuristic reasoning ([Goodfellow et al., 2016](#)). Once again, this principle is especially critical in sensitive domains where decisions must be accountable and ethically sound ([Mittelstadt, 2019](#); [Varshney and Alemzadeh, 2017](#)), and in NLP, where the inherent ambiguity and subjectivity of language complicate evaluation ([Hovy and Prabhu-moye, 2021](#)). In many NLP tasks (e.g., [Evkoski and Pollak, 2023](#)) a degree of approximation is often tolerated in favour of pragmatic utility, and models are evaluated based on what is useful or convincing to downstream consumers. By contrast, in our work, it is strongly desirable that model predictions remain traceable to the actual input provided by us. This grounding is not only central to scientific rigour, but also to ensuring justifiability and trust in use cases such as clinical assessments and the analysis of phenomenological interviews, where outputs may influence human understanding of complex experiences.

Importantly, ground-truth basis is complementary to interpretability. While interpretability focuses on making the model’s decision process understandable, ground-truth basis ensures that its outputs are substantively anchored in verified data rather than emergent patterns from opaque pretrain-

ing. Together, these two properties are essential for making computational predictions trustworthy and usable by stakeholders such as clinicians and phenomenologists.

As expected, achieving high classification **accuracy** remains a central objective, and considering all the other desiderata, a model with a lower **computational cost** is to be preferred. Additionally, given the sensitivity of the data, we prioritise tools that guarantee full control over processing and prevent third-party access.

7 Proposed Approaches (RO3b)

In this subsection, we refer to literature in order to compare the various proposed approaches with regard to each of our desiderata. The proposed approaches are: conventional discriminative models, including traditional AI and neural networks (NNs); generative LLMs, fine-tuned or with few-shot learning; and embedding-based retrieval approaches.

As the NLP landscape—particularly in relation to LLMs, interpretability, and domain-specific adaptation—continues to evolve rapidly, the methodological choices outlined below are intended as a flexible, revisable framework rather than a rigid pipeline. We anticipate that developments over the course of the Ph.D. will inform and potentially shift our implementation strategies, especially in response to emerging technologies and best practices in ethical, explainable NLP. In line with this flexible and modular approach, we also propose the investigation of a mixture-of-experts (MoE) architecture.

To train our models, we plan to employ both learned textual features—such as embeddings or TF-IDF representations—and predefined features derived from both previous studies (e.g., [Pennabaker et al., 2003](#)) and further investigations based on [Caporusso et al. \(2024\)](#)’s framework. This hybrid feature strategy supports both data-driven learning and interpretability through grounded linguistic markers.

Preliminary experiments are described in the Appendix A.

7.1 Conventional Discriminative Models

Conventional discriminative models include both traditional ML methods ([Bishop and Nasrabadi, 2006](#)) and NNs ([LeCun et al., 2015](#)). Examples include SVMs ([Cristianini and Shawe-Taylor, 2000](#)), Logistic Regression (LR), decision trees, and feed-

forward or recurrent NNs (RNNs) ([Goodfellow et al., 2016](#)) trained for classification purposes. They are often employed in the context of supervised learning, where the model learns from labelled data ([Murphy, 2012](#)).

Conventional discriminative models represent a good approach to our goal, assuming the availability of high-quality annotated datasets. Once trained, such models can directly classify a given text instance into predefined categories—such as BS, NS, or AS—and further specify the mode for each element (e.g., bodily ownership: present; agency over the body: partial). **Interpretability** in this approach depends largely on the choice of model: while rule-based models like decision trees or LR are inherently transparent, NNs are less interpretable and often require post-hoc explanation methods. Regarding **ground-truth** alignment, conventional discriminative models are optimal, since their outputs are entirely dependent on the patterns found in the labelled examples. When sufficient and representative training data is available, these models can be very **accurate**. Furthermore, they can be highly efficient **computationally**.

7.2 Generative LLMs

Generative LLMs (e.g., GPT; [Radford et al., 2018](#)) are designed to produce new outputs—in the case of language models, in the form of text—by learning the underlying distribution of the training data ([Bengio et al., 2003](#); [Radford et al., 2018](#)).

Although flexible, they come with a few challenges. For example, even in the case that their output looks plausible, it might be incorrect. This is referred to as *hallucination*, and it is due to the fact that these models generate responses solely based on learned statistical patterns ([Zhang et al., 2022](#)). Additionally, they reflect biases present in their training data and lack transparent mechanisms for interpreting or verifying their outputs ([Bolkunbasi et al., 2016](#)).

Ideally, generative LLMs will be applied to our task either through prompt-based few-shot learning or via fine-tuning on labelled datasets ([Wei et al., 2022](#); [Wolf et al., 2020](#)), which generally improves accuracy and control over outputs ([Howard and Ruder, 2018](#)).

While LLMs offer great flexibility and generalisation capabilities, they are not **interpretable**. Although post-hoc explanation methods like LIME (Local Interpretable Model-agnostic Explanations; [Alvarez-Melis and Jaakkola, 2018](#); [Ribeiro et al.,](#)

2016) or SHAP (SHapley Additive exPlanations; Jin et al., 2020; Lundberg and Lee, 2017) can provide some superficial insight, they do not guarantee true transparency or fidelity to the model’s internal reasoning. Furthermore, LLMs are not grounded in **ground-truth** data. Even when fine-tuned, it remains unclear whether these models’ predictions are derived from the data used for fine-tuning or the huge corpora used for pre-training. Furthermore, their outputs can change even from subtle shifts in prompt wording. This affects the consistency and reliability of the model. **Accuracy** is often high in LLMs, but it depends on prompt design and the complexity of the task. Inconsistent results could result from similar inputs, particularly when the classification schema is fine-grained, such as distinguishing between modes of Self-experience. Finally, generative LLMs are **computationally** expensive.

7.3 Embedding-Based Retrieval

Embedding-based retrieval is a type of retrieval-based approach which involves mapping the input into a shared vector space using models such as BERT (Devlin et al., 2019) or Sentence-BERT (Reimers and Gurevych, 2019). The vector representations of the inputs are compared to the already existing vector space, i.e., the knowledge base (Karpukhin et al., 2020). The initial vector space can be fine-tuned to task specific data, enhancing the model performance, and the semantic similarity between the reference and the input texts can be measured via cosine similarity or other distance metrics (Cer et al., 2018; Xiong et al., 2020).

For our purpose, embedding-based retrieval is especially useful in the case that a well-curated repository of annotated examples is available. The model can retrieve similar past instances that have already been labelled, allowing it to infer the classification of the new instance by analogy. While the embedding process itself is not inherently **interpretable**, the example-based reasoning enabled by retrieval models provides a form of implicit transparency: it is possible to inspect the retrieved examples and their labels to understand the basis of the model’s recommendation. This makes the approach more explainable than generative LLMs, although not as transparent as rule-based classifiers. In terms of **ground-truth** alignment, embedding-based retrieval performs strongly. The model’s decisions are anchored in annotated, verified data, and it does not generate new content but rather iden-

tifies the closest match among existing cases. In RAG-style architectures (retrieval-augmented generation; Lewis et al., 2020), this grounding helps reduce—but does not eliminate—the risk of hallucination during generation. **Accuracy** depends heavily on the quality and diversity of the dataset: if the database covers a broad range of expressions for different Self-aspects and modes, the model can achieve high classification performance. **Computationally**, this approach is efficient. Embeddings can be pre-computed, and retrieval operations (e.g., cosine similarity search) are lightweight.

7.4 Mixture of Experts

We also plan to explore a MoE architecture based on the work by Swamy et al. (2025), who proposed an interpretable MoE model designed for human-centric applications. In such architectures, different sub-networks—i.e., *experts*—are selectively activated depending on the input, enabling instance-specific reasoning and the possibility of **interpretability** where needed. This design offers a compelling balance between flexibility and transparency: it allows the integration of both interpretable models and black-box models within a unified framework. For our purposes, this means we can assign interpretable models to Self-aspect categories where explanation is critical (e.g., clinical applications), while using more complex models for noisier or less constrained categories.

The modular nature of MoE architectures also aligns well with our Self-aspect ontology. Since each expert can be specialised to a distinct subset of Self-aspects or linguistic patterns, this structure supports both conceptual clarity and efficient scalability (**computational cost**). Moreover, because only a few experts are activated per instance, the resulting predictions can offer local insight into the decision process, particularly when interpretable experts are selected. Importantly, expert modules trained on annotated data can maintain clear ties to their training supervision, preserving **ground-truth basis** at the module level. We believe this architecture is a promising direction to address the trade-off between **accuracy** and interpretability across the wide range of Self-related phenomena we aim to model.

8 Evaluation (RO4)

8.1 Intrinsic Evaluation

To evaluate and compare the effectiveness of different classification methods for identifying Self-aspects and their elements and modes in text, the approach proposed by [Demšar \(2006\)](#) to compare the performance of multiple classifiers across multiple datasets will be adopted. To use this method, a minimum of five different datasets is necessary, although it is recommended to employ at least 10. In the context of this Ph.D., a diverse range of models will be used to perform the classification (see Section 7). Despite their varied architectures and learning paradigms, they all can be evaluated in a comparable way. That is to say, by producing predictions over shared, annotated datasets and assessing them using standard performance metrics such as accuracy, F1-score, or macro-averaged precision and recall. By using [Demšar \(2006\)](#)'s framework, the evaluation will not only focus on raw performance, but also support robust conclusions about the relative strengths of each approach in the context of supervised Self-aspect classification. This is essential for making informed methodological choices, particularly when weighing the benefits of interpretable and ground-truth-aligned models against those of more flexible, data-driven generative LLMs. For the purposes of evaluation, we adopt an instance-based setup, treating each labelled unit (e.g., sentence or utterance) as a classification instance. Future work may explore span-based evaluation to capture finer-grained textual markers of Self-aspect expression.

8.2 Extrinsic Evaluation

We also plan to evaluate our framework by how useful it proves to be in downstream tasks. As it is likely that different trade-offs of desirable features are best for different applications, we do not aim to propose one singular model, but a collection of models. They will ideally be implemented in a user-friendly software that will allow the selection of the desired model, along with information and suggestions regarding each of them. Additionally, similarly to LIWC ([Boyd et al., 2022](#)), the user will be able to select which Self-aspects to analyse, and to which degree of granularity. It will be possible to determine at which level should the analysis be conducted, e.g., at the sentence, paragraph, or document level.

We plan to conduct at least two case studies in

which we will apply one or more of our developed models to different tasks.

In the context of an ongoing project on NLP approaches to cognitive decline, we plan to analyse comparable texts produced by clinical vs non-clinical population by using one or more of our proposed models. In particular, this will serve to test hypothesis on the differences in Self-aspects, but also, potentially, to identify features that could be used to detect cognitive decline.

In the context of the larger attempt to develop a computational framework to support the analysis of phenomenological interviews, one or more of our developed models will be adopted to support the analysis of the phenomenology of the Self, fundamental to most, if not all, experiences. This could help highlight how the Self is experienced differently across an episode (e.g., a dissolution experience; [Caporusso, 2022](#)), or how it is experienced by different populations, e.g., affected or not by derealisation.

9 Conclusion

We presented a proposal to design a computational model capable of detecting Self-aspects in text, grounded in a structured ontology and supported by diverse, annotated datasets curated by us. Our approach bridges conceptual insights from fields such as psychology and phenomenology with empirical techniques in NLP, enabling interpretable and application-oriented analysis of Self in language. Rather than relying on a single architecture, we propose and evaluate a range of computational models—rule-based, embedding-based, and generative LLMs—each assessed in light of desiderata such as interpretability, ground-truth basis, accuracy, and computational cost. By aligning technical development with ethical considerations and application-specific constraints, we aim to contribute not only a functional model, but also a thoughtful framework for the computational study of the Self.

10 Limitations

Our work presents various limitations. The Self-aspects specified in our ontology may be insufficient or suboptimal for the range of tasks we intend to address. Additionally, although our datasets are heterogeneous, this may still be insufficient for generalisability—particularly across cultural contexts where expressions of Self may vary signif-

icantly. The heterogeneity of the datasets, along with the flexible granularity of labelling units, may also introduce inconsistencies. Furthermore, many of the computational approaches we propose require substantial resources, including large volumes of annotated data. Moreover, there is a risk of overfitting to the specific theoretical assumptions embedded in our ontology, particularly if it privileges certain conceptions of the Self over others, potentially narrowing the interpretive scope of our models. Reconciling the need for interpretability and ground-truth adherence with high classification performance remains a central challenge in our methodological design.

11 Ethical Considerations

As this study relies on the use of existing datasets or datasets collected within the scope of other projects, the ethical considerations pertaining to each dataset are governed by the terms under which the data have been or will be collected. For datasets obtained through restricted access, we will comply with all necessary data use agreements and institutional requirements. We are committed to ensuring the anonymisation of all textual data prior to model training. Since our datasets and LLMs may reflect cultural or demographic biases, we acknowledge the risk of reproducing or amplifying such biases in our outputs. We emphasise that the computational models developed in this research are intended to function as support tools rather than as standalone decision-makers.

References

- Jonathan M Adler. 2012. Living into the story: agency and coherence in a longitudinal study of narrative identity development and mental health over the course of psychotherapy. *Journal of personality and social psychology*, 102(2):367.
- Jonathan M Adler, Lauren M Skalina, and Dan P McAdams. 2008. The narrative reconstruction of psychotherapy and psychological health. *Psychotherapy research*, 18(6):719–734.
- Muhammad Aurangzeb Ahmad, Carly Eckert, and Ankur Teredesai. 2018. Interpretable machine learning in healthcare. In *Proceedings of the 2018 ACM international conference on bioinformatics, computational biology, and health informatics*, pages 559–560.
- David Alvarez-Melis and Tommi S Jaakkola. 2018. On the robustness of interpretability methods. *arXiv preprint arXiv:1806.08049*.
- Julia Amann, Alessandro Blasimme, Effy Vayena, Dietmar Frey, Vince I Madai, and Precise4Q Consortium. 2020. Explainability for artificial intelligence in healthcare: a multidisciplinary perspective. *BMC medical informatics and decision making*, 20:1–9.
- Patrik Aspers. 2009. Empirical phenomenology: A qualitative research approach (the cologne seminars). *Indo-pacific journal of phenomenology*, 9(2).
- Michael Bamberg. 2008. Considering counter narratives. In *Considering counter-narratives: Narrating, resisting, making sense*, pages 351–371. John Benjamins Publishing Company.
- Roy F Baumeister, Arlene M Stillwell, and Todd F Heatherton. 1994. Guilt: an interpersonal approach. *Psychological bulletin*, 115(2):243.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155.
- José Luis Bermúdez. 2018. *The bodily self: Selected essays*. MIT Press.
- Christopher M Bishop and Nasser M Nasrabadi. 2006. *Pattern recognition and machine learning*, volume 4. Springer.
- Lasse Bohlen, Julian Rosenberger, Patrick Zschech, and Mathias Kraus. 2024. Leveraging interpretable machine learning in intensive care. *Annals of Operations Research*, pages 1–40.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.
- Ryan L Boyd, Ashwini Ashokkumar, Sarah Seraj, and James W Pennebaker. 2022. The development and psychometric properties of liwc-22. *Austin, TX: University of Texas at Austin*, 10:1–47.
- Ryan L Boyd and H Andrew Schwartz. 2021. Natural language analysis and the psychology of verbal behavior: The past, present, and future states of the field. *Journal of Language and Social Psychology*, 40(1):21–41.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Jerome Seymour Bruner. 2003. *Making stories: Law, literature, life*. Harvard University Press.
- Jaya Caporusso. 2022. Dissolution experiences and the experience of the self: an empirical phenomenological investigation (unpublished master’s thesis). university of vienna. *Advisor: Assist. Prof. Dr. Maja Smrdu*.

877	Jaya Caporusso, Boshko Koloski, Maša Rebernik,	Schwartz. 2018. Facebook language predicts depression in medical records. <i>Proceedings of the National Academy of Sciences</i> , 115(44):11203–11208.	931
878	Senja Pollak, and Matthew Purver. 2024. A		932
879	phenomenologically-inspired computational analysis		933
880	of self-categories in text. In <i>Proceedings of JADT</i>		
881	2024.		
882	Jaya Caporusso, Thi Hong Hanh Tran, and Senja Pollak.	Christopher Ifeanyi Eke, Azah Anir Norman, Liyana	934
883	2023. IJS@LT-EDI : Ensemble approaches to detect	Shuib, and Henry Friday Nweke. 2019. A survey	935
884	signs of depression from social media text . In <i>Pro-</i>	of user profiling: State-of-the-art, challenges, and	936
885	<i>ceedings of the Third Workshop on Language Tech-</i>	solutions. <i>IEEE Access</i> , 7:144907–144924.	937
886	<i>nology for Equality, Diversity and Inclusion</i> , pages		
887	172–178, Varna, Bulgaria. INCOMA Ltd., Shoumen,	Mohamad El Haj, Pascal Antoine, Jean Louis Nandrino,	938
888	Bulgaria.	and Dimitrios Kapogiannis. 2015. Autobiographical	939
889		memory decline in alzheimer’s disease, a theoret-	940
890	Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua,	ical and clinical overview. <i>Ageing research reviews</i> ,	941
891	Nicole Limtiaco, Rhomni St John, Noah Constant,	23:183–192.	942
892	Mario Guajardo-Cespedes, Steve Yuan, Chris Tar,		
893	et al. 2018. Universal sentence encoder. <i>arXiv</i>	Mohammad Ennab and Hamid Mcheick. 2024. En-	943
	<i>preprint arXiv:1803.11175</i> .	hancing interpretability and accuracy of ai models in	944
894		healthcare: a comprehensive review on challenges	945
895	Muhammad Waqas Anjum Ch and Waqas Arshad	and future directions. <i>Frontiers in Robotics and AI</i> ,	946
896	Cheema. 2018. A study of content based methods	11:1444763.	947
897	for author profiling in multiple genres. <i>International</i>		
898	<i>Journal of Scientific Engineering Research</i> , 9(9):322–	Bojan Evkoski and Senja Pollak. 2023. Xai in com-	948
	327.	putational linguistics: Understanding political lean-	949
899		ings in the slovenian parliament. <i>arXiv preprint</i>	950
900	Ahmad Chaddad, Jihao Peng, Jian Xu, and Ahmed	<i>arXiv:2305.04631</i> .	951
901	Bouridane. 2023. Survey of explainable ai tech-		
	niques in healthcare. <i>Sensors</i> , 23(2):634.	Mehwish Fatima, Komal Hasan, Saba Anwar, and Rao	952
902		Muhammad Adeel Nawab. 2017. Multilingual author	953
903	Mara Chineia-Rios, Thomas Müller, Gretel Liz De la	profiling on facebook. <i>Information Processing &</i>	954
904	Peña Sarracén, Francisco Rangel, and Marc Franco-	<i>Management</i> , 53(4):886–904.	955
905	Salvador. 2022. Zero and few-shot learning for au-		
906	thor profiling. In <i>International Conference on Appli-</i>	Mark Freeman. 2009. <i>Hindsight: The promise and peril</i>	956
907	<i>cations of Natural Language to Information Systems</i> ,	<i>of looking backward</i> . Oxford University Press.	957
	pages 333–344. Springer.		
908		Shaun Gallagher. 2000. Philosophical conceptions of	958
909	Nello Cristianini and John Shawe-Taylor. 2000. <i>An</i>	the self: implications for cognitive science. <i>Trends</i>	959
910	<i>introduction to support vector machines and other</i>	<i>in cognitive sciences</i> , 4(1):14–21.	960
911	<i>kernel-based learning methods</i> . Cambridge univer-		
	sity press.	Xiaoquan Gao, Sabriya Alam, Pengyi Shi, Franklin	961
912		Dexter, and Nan Kong. 2023. Interpretable machine	962
913	Janez Demšar. 2006. Statistical comparisons of clas-	learning models for hospital readmission prediction:	963
914	sifiers over multiple data sets. <i>Journal of Machine</i>	a two-step extracted regression tree approach. <i>BMC</i>	964
	<i>learning research</i> , 7(Jan):1–30.	<i>medical informatics and decision making</i> , 23(1):104.	965
915			
916	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Ian Goodfellow, Yoshua Bengio, Aaron Courville, and	966
917	Kristina Toutanova. 2019. Bert: Pre-training of deep	Yoshua Bengio. 2016. <i>Deep learning</i> , volume 1.	967
918	bidirectional transformers for language understand-	MIT press Cambridge.	968
919	ing. In <i>Proceedings of the 2019 conference of the</i>		
920	<i>North American chapter of the association for com-</i>	Sharath Chandra Guntuku, Rachelle Schneider, Arthur	969
921	<i>putational linguistics: human language technologies,</i>	Pelullo, Jami Young, Vivien Wong, Lyle Ungar,	970
	<i>volume 1 (long and short papers)</i> , pages 4171–4186.	Daniel Polsky, Kevin G Volpp, and Raina Merchant.	971
922		2019. Studying expressions of loneliness in individu-	972
923	Finale Doshi-Velez and Been Kim. 2017. Towards a	als using twitter: an observational study. <i>BMJ open</i> ,	973
924	rigorous science of interpretable machine learning.	9(11):e030355.	974
	<i>arXiv preprint arXiv:1702.08608</i> .		
925		Tilmann Habermas and Christin Köber. 2015. Auto-	975
926	Xiaowei Du and Yunmei Sun. 2022. Linguistic features	biographical reasoning in life narratives buffers the	976
927	and psychological states: A machine-learning based	effect of biographical disruptions on the sense of	977
	approach. <i>Frontiers in psychology</i> , 13:955850.	self-continuity. <i>Memory</i> , 23(5):664–674.	978
928			
929	Johannes C Eichstaedt, Robert J Smith, Raina M Mer-	Yaakov HaCohen-Kerner. 2022. Survey on profiling	979
930	chant, Lyle H Ungar, Patrick Crutchley, Daniel	age and gender of text authors. <i>Expert Systems with</i>	980
	Preoŕiuc-Pietro, David A Asch, and H Andrew	<i>Applications</i> , 199:117140.	981

982	Henriëtte Dorothée Heering, Saskia Goedhart, Richard	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	1034
983	Bruggeman, Wiepke Cahn, Lieuwe de Haan, René S	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	1035
984	Kahn, Carin J Meijer, Inez Myin-Germeys, Jim van	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	1036
985	Os, and Durk Wiersma. 2016. Disturbed experience	täschel, et al. 2020. Retrieval-augmented generation	1037
986	of self: psychometric analysis of the self-experience	for knowledge-intensive nlp tasks. <i>Advances in neu-</i>	1038
987	lifetime frequency scale (self). <i>Psychopathology</i> ,	<i>ral information processing systems</i> , 33:9459–9474.	1039
988	49(2):69–76.		
989	Tine Holm, Dorthe Kirkegaard Thomsen, and Vibeke	X Alice Li and Devi Parikh. 2019. Lemotif: An affec-	1040
990	Bliksted. 2016. Life story chapters and narrative	tive visual journal using deep neural networks. <i>arXiv</i>	1041
991	self-continuity in patients with schizophrenia. <i>Con-</i>	<i>preprint arXiv:1903.07766</i> .	1042
992	<i>sciousness and cognition</i> , 45:60–74.		
993	Dirk Hovy and Shrimai Prabhumoye. 2021. Five	Zachary C Lipton. 2018. The mythos of model inter-	1043
994	sources of bias in natural language processing. <i>Lan-</i>	pretability: In machine learning, the concept of in-	1044
995	<i>guage and linguistics compass</i> , 15(8):e12432.	terpretability is both important and slippery. <i>Queue</i> ,	1045
		16(3):31–57.	1046
996	Jeremy Howard and Sebastian Ruder. 2018. Universal	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	1047
997	language model fine-tuning for text classification.	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	1048
998	<i>arXiv preprint arXiv:1801.06146</i> .	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	1049
		Roberta: A robustly optimized bert pretraining ap-	1050
999	Baixiang Huang, Canyu Chen, and Kai Shu. 2024. Can	proach. <i>arXiv preprint arXiv:1907.11692</i> .	1051
1000	large language models identify authorship? <i>arXiv</i>		
1001	<i>preprint arXiv:2403.08213</i> .	Sheng-Chieh Lu, Christine L Swisher, Caroline Chung,	1052
		David Jaffray, and Chris Sidey-Gibbons. 2023. On	1053
1002	Baixiang Huang, Canyu Chen, and Kai Shu. 2025. Au-	the importance of interpretable machine learning pre-	1054
1003	thorship attribution in the era of llms: Problems,	dictions to inform clinical decision making in oncol-	1055
1004	methodologies, and challenges. <i>ACM SIGKDD Ex-</i>	ogy. <i>Frontiers in oncology</i> , 13:1129380.	1056
1005	<i>plorations Newsletter</i> , 26(2):21–43.		
1006	Clayton Hutto and Eric Gilbert. 2014. Vader: A pars-	Scott M Lundberg and Su-In Lee. 2017. A unified ap-	1057
1007	imonious rule-based model for sentiment analysis of	proach to interpreting model predictions. <i>Advances</i>	1058
1008	social media text. In <i>Proceedings of the international</i>	<i>in neural information processing systems</i> , 30.	1059
1009	<i>AAAI conference on web and social media</i> , volume 8,	Paul Henry Lysaker and John Timothy Lysaker. 2002.	1060
1010	pages 216–225.	Narrative structure in psychosis: Schizophrenia and	1061
		disruptions in the dialogical self. <i>Theory & Psychol-</i>	1062
1011	Jeff Jaeger, Katie M Lindblom, Kelly Parker-Guilbert,	ogy, 12(2):207–220.	1063
1012	and Lori A Zoellner. 2014. Trauma narratives: It’s		
1013	what you say, not how you say it. <i>Psychological</i>	Dan P McAdams. 2001. The psychology of life stories.	1064
1014	<i>Trauma: Theory, Research, Practice, and Policy</i> ,	<i>Review of general psychology</i> , 5(2):100–122.	1065
1015	6(5):473.		
1016	Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter	Seifeddine Mechti, Nabil Khoufi, and Lamia	1066
1017	Szolovits. 2020. Is bert really robust? a strong base-	Hadrich Belguith. 2020. Improving native language	1067
1018	line for natural language attack on text classification	identification model with syntactic features: Case	1068
1019	and entailment. In <i>Proceedings of the AAAI con-</i>	of arabic. In <i>Intelligent Systems Design and</i>	1069
1020	<i>ference on artificial intelligence</i> , volume 34, pages	<i>Applications: 18th International Conference on</i>	1070
1021	8018–8025.	<i>Intelligent Systems Design and Applications (ISDA</i>	1071
		<i>2018) held in Vellore, India, December 6-8, 2018,</i>	1072
1022	Vladimir Karpukhin, Barlas Oguz, Sewon Min,	<i>Volume 2</i> , pages 202–211. Springer.	1073
1023	Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi		
1024	Chen, and Wen-tau Yih. 2020. Dense passage re-	Marishka M Mehta, Soojung Na, Xiaosi Gu, James W	1074
1025	trieval for open-domain question answering. In	Murrough, and Laurel S Morris. 2023. Reward-	1075
1026	<i>EMNLP (1)</i> , pages 6769–6781.	related self-agency is disturbed in depression and	1076
		anxiety. <i>PloS one</i> , 18(3):e0282727.	1077
1027	Tine Kolenik, Günter Schiepek, and Matjaž Gams.	Matthias Michal, Bettina Reuchlein, Julia Adler, Iris	1078
1028	2024. Computational psychotherapy system for men-	Reiner, Manfred E Beutel, Claus Vögele, Hartmut	1079
1029	tal health prediction and behavior change with a	Schächinger, and Andre Schulz. 2014. Striking	1080
1030	conversational agent. <i>Neuropsychiatric Disease and</i>	discrepancy of anomalous body experiences with	1081
1031	<i>Treatment</i> , pages 2465–2498.	normal interoceptive accuracy in depersonalization-	1082
		derealization disorder. <i>PloS one</i> , 9(2):e89823.	1083
1032	Yann LeCun, Yoshua Bengio, and Geoffrey Hinton.		
1033	2015. Deep learning. <i>nature</i> , 521(7553):436–444.	Brent Mittelstadt. 2019. Principles alone cannot	1084
		guarantee ethical ai. <i>Nature machine intelligence</i> ,	1085
		1(11):501–507.	1086

1087	Hadi Mohammadi, Ayoub Bagheri, Anastasia Gi-	Claire Petitmengin, Anne Remillieux, and Camila	1139
1088	achanou, and Daniel L Oberski. 2025. Explainability	Valenzuela-Moguillansky. 2019. Discovering the	1140
1089	in practice: A survey of explainable nlp across vari-	structures of lived experience: Towards a micro-	1141
1090	ous domains. <i>arXiv preprint arXiv:2502.00837</i> .	phenomenological analysis method. <i>Phenomenology</i>	1142
		<i>and the Cognitive Sciences</i> , 18(4):691–730.	1143
1091	Christoph Molnar. 2020. <i>Interpretable machine learn-</i>	Valeria I Petkova, Malin Björnsdotter, Giovanni Gentile,	1144
1092	<i>ing</i> . Lulu. com.	Tomas Jonsson, Tie-Qiang Li, and H Henrik Ehrsson.	1145
1093	James W Moore. 2016. What is the sense of agency and	2011. From part-to whole-body ownership in the	1146
1094	why does it matter? <i>Frontiers in psychology</i> , 7:1272.	multisensory brain. <i>Current Biology</i> , 21(13):1118–	1147
1095	Kevin P Murphy. 2012. <i>Machine learning: a probabilis-</i>	1122.	1148
1096	<i>tic perspective</i> . MIT press.	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya	1149
1097	Jinkyung Na and Incheol Choi. 2009. Culture and first-	Sutskever, et al. 2018. Improving language under-	1150
1098	person pronouns. <i>Personality and Social Psychology</i>	standing by generative pre-training.	1151
1099	<i>Bulletin</i> , 35(11):1492–1499.	Matthew Ratcliffe. 2014. <i>Experiences of depression: A</i>	1152
1100	Ohad Nave, Fynn-Mathis Trautwein, Yochai Ataria,	<i>study in phenomenology</i> . OUP Oxford.	1153
1101	Yair Dor-Ziderman, Yoav Schweitzer, Stephen Fulder,	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	1154
1102	and Aviva Berkovich-Ohana. 2021. Self-boundary	Sentence embeddings using siamese bert-networks.	1155
1103	dissolution in meditation: A phenomenological in-	<i>arXiv preprint arXiv:1908.10084</i> .	1156
1104	vestigation. <i>Brain Sciences</i> , 11(6):819.	Marco Tulio Ribeiro, Sameer Singh, and Carlos	1157
1105	Matthew M Nour, Lisa Evans, David Nutt, and	Guestrin. 2016. Why should i trust you? explaining	1158
1106	Robin L Carhart-Harris. 2016. Ego-dissolution and	the predictions of any classifier. In <i>Proceedings of</i>	1159
1107	psychedelics: validation of the ego-dissolution in-	<i>the 22nd ACM SIGKDD international conference on</i>	1160
1108	ventory (edi). <i>Frontiers in human neuroscience</i> ,	<i>knowledge discovery and data mining</i> , pages 1135–	1161
1109	10:190474.	1144.	1162
1110	Sarra Ouni, Fethi Fkih, and Mohamed Nazih Omri.	Stephanie Rude, Eva-Maria Gortner, and James Pen-	1163
1111	2021. Toward a new approach to author profiling	nebaker. 2004. Language use of depressed and	1164
1112	based on the extraction of statistical features. <i>Social</i>	depression-vulnerable college students. <i>Cognition &</i>	1165
1113	<i>Network Analysis and Mining</i> , 11(1):59.	<i>Emotion</i> , 18(8):1121–1133.	1166
1114	Sarra Ouni, Fethi Fkih, and Mohamed Nazih Omri.	Cynthia Rudin. 2019. Stop explaining black box ma-	1167
1115	2023a. Novel semantic and statistic features-based	chine learning models for high stakes decisions and	1168
1116	author profiling approach. <i>Journal of Ambient In-</i>	use interpretable models instead. <i>Nature machine</i>	1169
1117	<i>telligence and Humanized Computing</i> , 14(9):12807–	<i>intelligence</i> , 1(5):206–215.	1170
1118	12823.	Rupsa Saha, Ole-Christoffer Granmo, and Morten Good-	1171
1119	Sarra Ouni, Fethi Fkih, and Mohamed Nazih Omri.	win. 2023. Using tsetlin machine to discover inter-	1172
1120	2023b. A survey of machine learning-based au-	pretable rules in natural language processing applica-	1173
1121	thor profiling from texts analysis in social networks.	tions. <i>Expert Systems</i> , 40(4):e12873.	1174
1122	<i>Multimedia Tools and Applications</i> , 82(24):36653–	Rupsa Saha, Ole-Christoffer Granmo, Vladimir I	1175
1123	36686.	Zadorozhny, and Morten Goodwin. 2022. A rela-	1176
1124	Josef Parnas, Paul Møller, Tilo Kircher, Jørgen Thal-	tional tsetlin machine with applications to natural	1177
1125	bitzer, Lennart Jansson, Peter Handest, and Dan Za-	language understanding. <i>Journal of Intelligent Infor-</i>	1178
1126	havi. 2005. Ease: examination of anomalous self-	<i>mation Systems</i> , pages 1–28.	1179
1127	experience. <i>Psychopathology</i> , 38(5):236.	H Andrew Schwartz, Johannes C Eichstaedt, Mar-	1180
1128	James W Pennebaker and Sandra K Beall. 1986. Con-	garet L Kern, Lukasz Dziurzynski, Stephanie M Ra-	1181
1129	fronting a traumatic event: toward an understanding	mones, Megha Agrawal, Achal Shah, Michal Kosin-	1182
1130	of inhibition and disease. <i>Journal of abnormal psy-</i>	ski, David Stillwell, Martin EP Seligman, et al. 2013.	1183
1131	<i>chology</i> , 95(3):274.	Personality, gender, and age in the language of social	1184
1132	James W Pennebaker and Martha E Francis. 1996. Cog-	media: The open-vocabulary approach. <i>PloS one</i> ,	1185
1133	gnitive, emotional, and language processes in disclo-	8(9):e73791.	1186
1134	sure. <i>Cognition & emotion</i> , 10(6):601–626.	Andrea Serino, Adrian Alsmith, Marcello Costan-	1187
1135	James W Pennebaker, Matthias R Mehl, and Kate G	tini, Alisa Mandrigin, Ana Tajadura-Jimenez, and	1188
1136	Niederhoffer. 2003. Psychological aspects of natural	Christophe Lopez. 2013. Bodily ownership and self-	1189
1137	language use: Our words, our selves. <i>Annual review</i>	location: components of bodily self-consciousness.	1190
1138	<i>of psychology</i> , 54(1):547–577.	<i>Consciousness and cognition</i> , 22(4):1239–1252.	1191

phenomenological expertise. Inter-annotator agreement was assessed via Cohen’s Kappa: 0.80 between human annotators, and 0.84–0.89 between human and model annotators.

A.2 Experimental Setup

We trained and evaluated six models using 10-fold cross-validation, combining three different classifiers—SVM, Logistic Regression (LR), and Naïve Bayes (NB)—with two types of feature representations. The first type comprised learned features, specifically TF-IDF weighted unigrams and bigrams. The second relied on predefined features derived from the LIWC-22 lexicon, specifically those previously identified as correlated with the Social Self aspect (Caporusso et al., 2024). Text preprocessing included converting all text to lowercase, removing punctuation, and applying z-score normalisation to the LIWC-derived features to ensure comparability across feature scales. To interpret the trained models, we employed feature importance techniques tailored to each algorithm: linear SVM coefficients for SVM, SHAP values for Logistic Regression, and permutation importance for Naïve Bayes.

A.3 Results

The best-performing model was the SVM trained on LIWC features, achieving a precision of 0.81 (STD = 0.03), recall of 0.82 (STD = 0.02), and F1-score of 0.81 (STD = 0.03) across 10 folds. It consistently outperformed all other models. Models using learned features (TF-IDF) performed slightly worse overall, with the SVM on learned features achieving an F1-score of 0.73 (STD = 0.04) and particularly lower recall. Statistical analysis confirmed the significance of these differences via a Friedman test (statistic = 44.26, $p < 0.001$) and pairwise Wilcoxon signed-rank tests (adjusted $p = 0.03$ for several comparisons). Feature importance analyses identified intuitive and interpretable markers of Social Self, including "we", social referents, affect terms, and pronoun use, aligning with prior findings and theoretical expectations.

A.4 Implications and Limitations

This pilot study demonstrates that interpretable models trained on psychologically grounded features can reliably identify expressions of Social Self in everyday texts. It also confirms the utility of a hybrid human-LLM annotation pipeline,

especially in early dataset development. However, several limitations emerged. Performance is currently limited to binary classification of a single Self-aspect. The current study also relies on English-language data, which restricts immediate generalisability.