

A computational approach to typological comparative concepts for lexicality

Coleman Haley



Doctor of Philosophy

Institute for Language, Cognition and Computation

School of Informatics

University of Edinburgh

2025

Abstract

One major dimension of linguistic organization is the notion that there are more lexical linguistic units, which express meanings, and more functional linguistic units, which are determined by syntax and/or discourse and serve to organize and clarify the relationships between lexical elements. This dichotomy has been described at levels of linguistic structure and motivates at least two classical distinctions in linguistics. At the level of words, it motivates the so-called lexical-functional distinction, while within morphology, a related distinction is drawn between derivation (which forms new lexical items) and inflection (which produces forms of lexical items). These dichotomies have many noted boundary cases, which have led to many linguists rejecting them, or treating them as gradient. In this thesis, I refer to this gradient of semantic weight at different levels of formal structure as lexicality.

There is substantial neurological and psychological evidence for the importance of lexicality to human language processing. Further, lexicality dichotomies also emerge in cross-linguistic trends in grammatical organization, such as asymmetries between inflection and derivation, or between the properties of functional and lexical word classes. Yet the lexicality of a particular linguistic unit varies contextually and diachronically. I develop quantitative methods to test the consistency of these concepts across typologically diverse languages. First, I show inflection vs. derivation can be predicted with high accuracy from formal and distributional properties.

In linguistic practices that proceed from analysis of language-particular data to a language-general analysis, issues of lexicality have played a role of central importance. However, in the functional—typological tradition, which proceeds from cross-linguistic analysis to the language particular, the relationship of this dimension to linguistic organization has had little theoretical impact. A major factor is that typological research must be conducted with cross-linguistically

applicable comparative concepts. In this thesis, I leverage deep learning models to produce empirically grounded measures for lexicality, which I argue can serve as interesting and useful comparative concepts for typological study.

In the first part of the thesis, I focus on inflection and derivation, operationalizing a four-dimensional framework for formal and distributional properties of the distinction. I show that formal and distributional variability are strong correlates of this traditional distinction across a sample of 26 languages, and that the four measures can predict inflection vs. derivation with 90% accuracy

In the second part of the thesis, I introduce a novel groundedness measure, which aims to provide a cross-linguistic empirical ground for language function to quantify contextual semantic contentfulness. To do so, I leverage image-caption datasets and vision-language models. This measure captures the lexical-functional distinction in word classes across 30 languages but diverges substantially from related measures like concreteness.

Interestingly, groundedness displays asymmetries not just between lexical and functional items, but also among the major lexical classes of nouns, verbs, and adjectives. I argue that this suggests a connection between ideas of lexical word-class continua in cognitive linguistics and the lexical-functional distinction. I apply groundedness to deviations from prototypical lexical class organization. I show that groundedness predicts the split between Japanese *na*- and *i*-adjectives, which has previously been thought to have little synchronic relevance. On the other hand, an investigation of the Tensedness Hypothesis shows the challenges with certain types of cross-linguistic comparisons of groundedness with current methods.

Lay Summary

Lay summary here

Acknowledgements

Acknowledgements here

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Coleman Haley)

Contents

List of Figures	xiii
------------------------	-------------

List of Tables	xvii
-----------------------	-------------

1 Introduction	1
1.1 The role of lexicality in linguistic organization	1
1.2 Approach	4
1.3 Structure of the Thesis	6
1.4 Contributions	11
2 Background	13
2.1 Typology and Comparative Concepts	14
2.1.1 Defining Hybrid Comparative Concepts	17
2.2 Finding Meaning in Computational Models	23
2.2.1 Distributional Embeddings	24
2.2.2 Contextual Embeddings and Language Models	24
2.2.3 Multimodal Models	24
2.3 Approaches to Comparative Concepts in Computational Typology	25
2.3.1 Comparative concepts in multilingual databases	25
2.3.2 Phonological typology	26
2.3.3 Semantic category systems	28
2.3.4 Multidimensional scaling	29

2.3.5	Deep learning models of comparative concepts	32
2.4	Formal and functional dimensions of lexicality	34
2.4.1	The formal dimension	35
2.4.2	The functional-semantic dimension	38
2.4.3	Summary	42
2.5	Chapter Summary	43
I	Inflection and Derivation	45
3	Corpus-based Measures for Inflection and Derivation	47
3.1	Introduction	47
3.2	Motivation for our measures	52
3.3	Method	56
3.3.1	Orthography-based measures	57
3.3.2	Distributional-embedding-based measures	58
3.4	Data	62
3.4.1	Data selection and summary	64
3.5	Distribution of the individual measures	66
3.5.1	Effects of Frequency	68
3.6	The role of syntactic information	69
4	Predicting Inflection and Derivation Cross-linguistically	73
4.1	Predicting inflection and derivation	73
4.2	Classification of Linguistic Types of Inflection	77
4.2.1	Categories of inflectional meaning	78
4.2.2	Inherent vs. contextual inflection and transpositions . .	81
4.2.3	Summary	83
4.3	Discussion	84
4.3.1	The role of our individual measures	84

4.3.2	Language generality	87
4.3.3	The classification approach	89
4.3.4	Inflection and derivation: gradient or categorical?	91
4.3.5	Are inflection and derivation identifiable from the statistics of language?	92
4.3.6	Classification and syntactic change	95
4.3.7	Future work	96
4.4	Conclusion	98
II	Word Classes	101
5	Groundedness and the Lexical–Functional Distinction	103
5.1	Introduction	103
5.2	Background	106
5.2.1	Contentfulness and word class	107
5.2.2	Measuring contentfulness	107
5.3	Method	109
5.4	Experimental setup	111
5.5	Results	114
5.5.1	Which word classes are grounded?	114
5.5.2	Which word classes are more grounded?	116
5.5.3	How consistent is word class groundedness across languages?	118
5.5.4	Semantic dimension of the measure	120
5.6	Discussion and Conclusion	122
6	Splitting and lumping: Visual groundedness as an organizing factor among lexical classes	127
6.1	Introduction	127

6.2	Continua among lexical word classes	130
6.3	Japanese adjectives	134
6.3.1	Method	136
6.3.2	Results	138
6.3.3	Discussion	140
6.4	The Tensedness Correlation	143
6.4.1	The typological finding	144
6.4.2	Theoretical explanation of the finding	146
6.4.3	Methodological background	147
6.4.4	Results	154
6.4.5	Discussion	155
6.5	Conclusion	159
	Bibliography	161

List of Figures


2.1	An example semantic map for the dative domain, adapted from Haspelmath (2003). Nodes represent different functions which “dative-like” elements can express. The boundaries for English <i>to</i> and French <i>à</i> are shown in pink and blue, respectively. Both terms cover contiguous regions of the map, satisfying the Semantic Map Connectivity Hypothesis.	19
3.1	The empirical distributions of our four measures (quantifying the magnitude <i>M</i> and variability <i>V</i> of changes in Form and in Embedding space) for inflections and derivations in UniMorph	67
3.2	The mean cosine similarity between FastText embeddings of words of the same and different parts of speech in UniMorph. .	71
4.1	Cross-validation accuracy and standard error in reconstructing UniMorph’s inflection–derivation distinction by various supervised classifiers. Linguistically-motivated hypotheses referred to in the text are denoted with letters	76
4.2	Probability and Odds ratio with 95% confidence intervals of being classified as derivation for various kinds of inflectional meaning. Inflections to the right of the dotted line were disproportionately likely to be classified as derivation by our model	79

4.3	Probability and Odds ratio with 95% confidence intervals of being classified as derivation for inherent inflections and transpositions	81
4.4	Probability and Odds ratio with 95% confidence intervals of being classified as derivation for inherent vs. contextual noun inflections	83
4.5	Our two most predictive measures for inflection and derivation. Saturation represents overlapping constructions. With respect to these two variables, the inflection–derivation distinction appears gradient rather than categorical	93
5.1	Mean and standard deviation of per-language mutual information estimates between word class and image. Across 30 languages, we see clear and consistent tendencies about which parts of speech are more “grounded”, corresponding to a distinction between lexical and functional classes.	104
5.2	Heatmap of mutual information estimates across parts of speech in thirty languages. Cells show the statistical significance of a word class’s groundedness ($MI > 0$). Unattested classes are white. Some functional classes display non-significant levels of groundedness in several languages, while lexical classes dominantly show highly significant grounding.	115
5.3	Word token level distributions of the groundedness measure (PMI) across all languages and datasets, grouped by part of speech (word class). We also report the estimated marginal mean and ranking of each word class. Colors are based on the ranking of classes, rather than their average PMIs. Overall, the distribution and estimated ranking of word classes strongly suggest our groundedness measure quantitatively captures the distinction between lexical and functional classes.	117

5.4	Correlation between human concreteness ratings and type-level groundedness (PMI; left, $\rho = 0.368$) or uncertainty coefficient (right, $\rho = 0.609$): i.e., the average ratio between LM surprisal and captioning model surprisal.	119
6.1	Groundedness scores for <i>na</i> -adjective <i>makka</i> (completely red; right) and <i>i</i> -adjective <i>akai</i> (red; left) in the STAIR-full-dev dataset.	142
6.2	Groundedness of the verbal categories across the 30 languages in this study. Error bars represent standard error in the mean groundedness across the datasets considered (COCO-35L, XM3600, and Multi30K where available). Contra theoretical predictions, verby languages do not exhibit higher mean groundedness of verbs, but are somewhat below average. However, this effect is confounded by model quality issues, as suggested by the lower groundedness of verbs in non-Latin script languages.	153
6.3	Z-scored groundedness of the verbal categories. Error bars represent standard error in the Z-scores across the datasets considered (COCO-35L, XM3600, and Multi30K where available). The results suggest verbs are not <i>relatively</i> more grounded than other words in verby languages. However, we observe a clear effect of script, with languages written in Latin script exhibiting relatively more grounded verbs.	156

6.4 Z-scored groundedness of the verbal categories, with adjectives included for verby languages. Error bars represent standard error in the Z-scores across the datasets considered (COCO-35L, XM3600, and Multi30K where available). Despite the higher groundedness of adjectives than verbs in general, and concerns that legitimate members of the verbal category could be disproportionately “lost” to the adjective tag in verby languages, we still observe lower groundedness for the verby languages. This suggests an disproportionate effect of captioning and language model quality on verbs. 157

List of Tables

3.1	Sample of an inflectional construction (upper table, German nominative plural) and derivational construction (lower table, English verbal nominalisation with <i>-ion</i>) in our data	55
3.2	Descriptive statistics of our filtered dataset by language.	65
5.1	We match the data points on which the language model and image captioning model were trained. The three datasets are the Gemma pre-training mixture (PT), PaliGemma multimodal data for continued training (CT), and COCO image–caption pairs for fine-tuning (FT). Symbols indicate whether models are trained on text data (A) or on multimodal data ( A).	111
6.1	?’s analysis of the conceptual categories of the major parts of speech and their semantic properties.	131
6.2	Differences in groundedness between adjective classes across datasets. “MT?” indicates whether the captions were machine-translated from English. The effect size is the increase in groundedness (in bits) associated with <i>na</i> -adjective-hood, estimated using a linear mixed effects model with fixed effects of word class and position and a random effect for word type. Overall, <i>na</i> -adjectives tend to be more grounded than <i>i</i> -adjectives. (<u>Significant results</u>)	139

6.3	The effect of adjective class on LM surprisal and captioning surprisal. We find that <i>na</i> -adjectives tend to be more surprising in the language model than <i>i</i> -adjectives, but this effect is reduced by conditioning on the images, resulting in higher overall groundedness. (<u>Significant results</u>)	141
-----	---	-----

Chapter 1

Introduction

1.1 The role of lexicality in linguistic organization

The distinction between LEXICAL and FUNCTIONAL linguistic units has played a role in the theory and analysis of language for millennia. This is to say, a distinction can be drawn between two poles for the role that linguistic units play in communication. On one end, we have the LEXICAL: linguistic signs which carry specific meanings, often referring out to objects or events in the world—nouns like *cat* or *tree*. On the other end, we have the FUNCTIONAL: signs which do not so much carry specific meanings, but rather serve to organize and clarify the relationships between lexical elements—like the tense marker *-ing* or the word *to*.

Because the roles served by FUNCTIONAL units are similar to roles that, in some languages, are expressed not through independent units but through structural patterns or rules (that is, through *grammar*), such units are sometimes referred to as GRAMMATICAL units. For example, in English, whether a word serves as the subject or object of a verb is indicated through word order alone, while in some languages, there are linguistic signs (called, variously *case markers*, *adpositions*, or *flags*) which explicitly mark these relationships. As such, for

theories which treat grammar as a separate system from the lexicon, functional units present a key challenge, as they straddle the boundary between these two systems, having the realized form of a linguistic sign, but the organizational role of grammar.

However, while there is some consistency in what concepts can be expressed as obligatory, paradigmatic, bound markings across languages, there are also serious definitional issues around where the boundary between lexical and functional units lies. To claim that “only certain concepts can be expressed functionally” presupposes a consistent definition of what it means to be functional; for this to avoid circularity, the definition of functional must not rely on the concepts themselves, but rather on something about linguistic *distribution*.

Yet the formal and distributional properties of functional expression are far from clear-cut. Typically, these are identified as (I) being CLOSED-CLASS, (II) being BOUND, and (III) being OBLIGATORY. While functional categories are typically closed-class (resisting new members), prototypically lexical categories can also be closed class, like Bemba adjectives (Dixon, 1977) or Jaminjung verbs (Pawley, 2006). Further, closedness is not a binary property, with closed classes varying substantially both in their size and their resistance to admitting new members. Purported functional elements can vary significantly in their degree of boundness, and indeed “boundness” is a complex property of dubious categorical status (??), with no consensus on how to define or measure it. Finally, obligatoriness is also gradient, with some functional elements being optional in some or even many contexts. The problem grows only more complex when we consider that the lexical status of a linguistic unit can vary contextually and diachronically.

While typical discussions of the distinction between lexical and functional units tend to focus on a contrast between lexical words or roots and functional words or affixes, a closely related distinction is drawn *within* the domain of

morphology between DERIVATION and INFLECTION. Morphology described as *inflectional* typically have all the prototypical properties of functional units: inflection is typically closed-class, bound, and obligatory, and align closely with “possible grammatical concepts.” In contrast, *derivational* morphology may express rich and perhaps unconstrained meaning, interacting ideosyncratically with the meaning of roots, and is typically optional. However, there are a few key differences with the lexical–functional distinction at the level of words. First, there are strong cross-linguistic tendencies for derivations to occur closer to the root than inflections do (Greenberg’s Universal 28) (Greenberg, 1966a; ?). Second, simple conversions or transpositions of word class (e.g. converting an adjective to a noun: *happy*→*happiness*) tend to pattern more similarly to derivations than inflection (e.g. scoping inside inflections), despite their apparent lack of semantic weight and highly obligatory and productive nature. This has led to substantial debate over whether such morphological constructions are better considered inflection (?) or derivation (ten Hacken, 1994). Further, where exactly the boundary between “transpositions of word class” and more semantically rich derivations lies is also unclear. For example, the *-er* nominalizer in English forms agentive nouns from verbs (*teach*→*teacher*). Again, we have encountered a distinction between two poles that seem to be related to semantic weight, but with unclear boundaries. While derivational morphology tends to be more semantically rich than inflectional morphology, it is typically less semantically rich than lexical roots; however, in highly agglutinative languages like Inuktitut, derivational morphemes can carry semantic content comparable to roots in other languages.

This thesis concerns itself with these divisions between more lexical and more functional linguistic units, at multiple levels of linguistic structure. I treat both the lexical–functional distinction at the level of words and the inflection–derivation distinction at the level of morphology, and connect them to prototype

phenomena among the major lexical classes of nouns, verbs, and adjectives. *In this thesis, I refer to this gradient of semantic weight at different levels of linguistic structure as the spectrum of LEXICALITY.*

Despite a rich base of evidence for differential representation and processing across the lexicality spectrum in psycholinguistics and neurolinguistics (Laudanna et al., 1992; Kirkici and Clahsen, 2013; ?; ?; ?; ?; ?; ?), the direct study of semantic weight/force/contentfulness on linguistic structure remains largely pre-theoretical in linguistics, especially in large-scale cross-linguistic study. A major cause of this theoretical lacuna is the difficulty of specifying *semantic contentfulness* in a principled, cross-linguistically applicable way. This has led many linguists to avoid this notion entirely, focusing instead on how notions like frequency shape grammatical expression (?). In this thesis, I seek to address this gap by developing measures of the *semantic* dimensions of lexicality distinctions and investigating *how they relate* to traditional “problematic” cross-linguistic grammatical distinctions.

1.2 Approach

Multiple Levels of Linguistic Structure This thesis spans a *wide range* of levels of linguistic structure. While previous work has largely treated these distinctions at different formal and semantic levels as different, unrelated problems (e.g. the inflection–derivation distinction at the morphological level; the distinction between lexical and functional word classes at the word level; or the distinctions between the major lexical classes of nouns, verbs, and adjectives), I investigate lexicality across multiple levels of linguistic structure, showing new parallels and connections between them. That being said, I limit my focus to *sub-phrasal* linguistic units (morphemes and words), leaving phrases, semantic frames, and more schematic constructions to future work.

Cross-linguistic investigation This thesis focuses on the *cross-linguistic* consistency of lexicality distinctions.¹ A large body of work in linguistic typology has argued that language-specific categories do not map onto some clean set of universal grammatical categories (Haspelmath, 2007; Croft, 2001; Dixon, 1977). Instead, typologists have argued for the importance of cross-linguistically valid *comparative concepts*—which need not necessarily map onto the structure of individual language’s grammar (Haspelmath, 2010; Croft, 2016a). Studies that focus on the distinction between inflection and derivation or between lexical and functional word classes which consider only a single language risk conflating language-particular properties and categories with cross-linguistic generalizations. While finding a consistent distinction between inflection and derivation, or between lexical and functional word classes in an individual language is interesting, it does *not* a-priori tell us whether such a distinction has cross-linguistic descriptive value. Thus, this thesis aims to cover a large and diverse sample of languages wherever possible.

Computational and quantitative methods To study the cross-linguistic consistency of lexicality distinctions, I take inspiration from the successes of empirical grounding in certain areas of typology (like vowel, color, and kinship systems; ?; ?; ?) and from recent advances in deep learning models of language. These models have been shown to learn rich representations of semantics and the world, without requiring direct instruction on this structure, but rather learning it implicitly from learning to predict words in a (linguistic and/or visual) context. This capacity makes these tools ideal for operationalizing semantic dimensions of lexicality distinctions. Further, in the second part of the thesis I leverage *multimodal* models which ground languages in images. This image

¹While in Chapter 6, I do conduct a language-particular analysis of Japanese word classes, the motivation for this analysis is to investigate whether cross-linguistically validated measures of lexicality can explain unusual language-particular patterns.

grounding provides a language- and form-neutral representation of semantics, which enables a new method for separating out contextual contentfulness from formal linguistic predictability. The computational approach also aids in our goal of cross-linguistic investigation—while psychological and neurological evidence for lexicality distinctions exists, scaling it to typological study is challenging. Computational methods require only corpus data, which enables the simultaneous study of many languages, though it biases the study towards languages with sufficient digital resources. Further, the quantitative focus of this thesis enables the study and quantification of consistency, in contrast with many previous studies that focus primarily on problematic cases.

1.3 Structure of the Thesis

I investigate three key research questions in this thesis:

- RQ1: What is the interplay between *form* and *semantics* across the lexicality spectrum? (Chapters 4, 5, 6)
- RQ2: How can we operationalize semantic contentfulness in a cross-linguistically applicable way? (Chapters 3, 5)
- RQ3: How *cross-linguistically consistent* are lexicality-related divisions like the division between the lexical and functional word classes, or the inflection–derivation distinction? (Chapters 4, 5)

Chapter 2: Background In this chapter, I expand on the theoretical framework for the thesis. I introduce in more detail the problems of cross-linguistic category comparison, and the method of comparative concepts for resolving these issues in typological research. I review the ways in which semantic function has been handled in typological comparative concepts, highlighting a role for deep

learning models of language in the creation of functional² comparative concepts. I then provide an overview history of the ways comparative concepts have been (at times, implicitly) employed, handled, and studied in computational typology. Through this background, I highlight how the creation of empirically grounded comparative concepts through the development of technologies and measures outside of typology (like perceptual theories of vowels and colors) has enabled major advances in typological research in the domains where it has been possible. This serves as further motivation for the computational approach to comparative concepts taken in this thesis. Finally, I provide a broad-scale overview of the lexicality spectrum, identifying different manifestations of a correlation between semantic contentfulness and formal linguistic structure, providing the connective tissue for the studies that follow. More detailed background on specific lexicality-related distinctions is provided in the relevant chapters.

Part I: Inflection and Derivation

Part I of this thesis consists of Chapters 3–4, which focus on the inflection–derivation distinction drawn in morphology. These chapters are based primarily on the following journal article:

Haley, C., Ponti, E. M., and Goldwater, S. (2024). Corpus-based measures discriminate inflection and derivation cross-linguistically. *Journal of Language Modelling*, 12(2):477–529

Chapter 3: Corpus-based Measures for Inflection and Derivation In this chapter, we introduce a computational framework for the inflection–derivation distinction. Inspired by Spencer (2013)’s description of the distinction, we introduce a set of four quantitative measures of morphological constructions, including measures of both the magnitude and the variability of the changes to

²In this sense, semantic

form and *distribution* introduced by each construction. Crucially, these measures can be computed directly from a linguistic corpus, allowing us to consistently operationalise them across many languages and morphological constructions. In contrast to prior computational studies that focus on a single language, we investigate 26 languages using the UniMorph 4.0 corpus (Batsuren et al., 2022). Using these measures, we find differences between inflection and derivation for all four measures, but substantial overlap for each individual measure. We demonstrate that the measures are not explained by simple frequency effects, and that the distributional measures capture a limited amount of syntactic information in addition to semantic information.

Chapter 4: Predicting Inflection and Derivation Using the measures from Chapter 3, we train classifier models to predict whether a construction is labeled as inflection or derivation in UniMorph. We find that language-agnostic classifier models over our measures are able to predict inflectional–derivational status with high accuracy (90%). We investigate linguistic categories of inflection, finding inflectional transpositions like participles are *not* more likely to be misclassified as derivational, in line with ?’s argument that these are best considered inflectional. Overall, our results are in line with a *consistent, yet gradient* view of the inflection–derivation distinction. Our results suggest that distributional and formal *variability* are the most important dimensions for the distinction, but the magnitude of distributional and formal change also play a role. While there is substantial overlap between the two categories on each individual dimension, the combination of all four dimensions provides a robust signal for the distinction in our sample.

Tricky knot to tie–magnitude seems less important, which is kind of at odds with the intro.

Part II: Word Classes

Part II of this thesis consists of Chapters 5–6, which focus on lexicality among (functional *and* lexical) word classes.

Chapter 5: Groundedness and the Lexical-Functional Distinction In this chapter, I introduce *groundedness*, a new semantic-contentfulness measure based on multimodal models. Focusing on the domain of image captions, I am able to treat an image as a proxy for a caption’s meaning. Using a language model and an image captioning model, I am able to estimate the pointwise mutual information between a token and the image as a surprisal difference under the two models. In this chapter, I focus on the **lexical-functional distinction** in parts of speech.

Using image captioning data in 30 languages from 10 language families, I find this groundedness measure largely rediscovers the distinction between lexical and functional word classes across 30 languages. Further, though it correlates only weakly with norms like imageability and concreteness in English, it provides a ranking suggested by cognitive linguists between nouns, verbs, and adjectives (noun > adjectives > verbs) across languages but contradicts the view of adpositions as a “semi-lexical” class. However, our results suggest grammatical word classes still carry semantic content. These results suggest the utility of this measure as a general tool for studying contentfulness in linguistics, and of taking a visually grounded approach to typological problems. This chapter is based on a conference paper at which appeared at NAACL 2025:

Haley, C., Goldwater, S., and Ponti, E. M. (2025). A Grounded Typology of Word Classes. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 10380–10399, Albuquerque, New Mexico. Association for Computational Linguistics.

Chapter 6: Splitting and Lumping In this chapter, I investigate the relationship between visual groundedness and cross-linguistic variation in **lexical parts of speech**. While there has been substantial work in linguistic typology investigating cross-linguistic variation in the expression of major lexical categories like nouns, verbs, and adjectives, this work has previously been largely disconnected from work on semantic contentfulness, the lexical–functional distinction, and grammaticalisation. Building on the visual groundedness measure introduced in Chapter 5, I connect existing continuum and prototype theories of lexical categories and meanings with groundedness. I argue that the role of semantic contentfulness in lexical categories can help explain cross-linguistic variation in lexical category organisation. In particular, I focus on languages which have been argued to “split” or “lump” major lexical categories.

To establish this, I first investigate Japanese. In Japanese, words denoting “properties” have the unusual property of constituting two formally very distinct word classes, rather than a single “adjective” class. Building on the insight that one of these classes is more formally “nominal” (*na*-adjectives) and one more “verbal” (*i*-adjectives), I hypothesise that we should see analogous trends in function: one class serving more prototypically nominal functions and one more prototypically verbal. In terms of visual groundedness, this corresponds to higher values for the nominal class. I investigate two manually captioned datasets and one machine translated dataset, finding significantly higher groundedness for *na*-adjectives in the manually captioned datasets, in line with the theoretical predictions. This stands in contrast to previous studies, which have indicated little synchronic functional difference between the two classes.

To investigate lumping phenomena, I turned to the Tensedness Correlation, which correlates the formal similarities of adjectives to verbs in languages with a lack of obligatory tense marking on verbs. In languages with obligatory tense marking, the expression of adjectives is more similar to nouns. I investigate

whether this correlation is reflected in groundedness, drawing on previous hypotheses for the cause of the Correlation. I find no significant relationship, which I argue is due to issues with directly comparing groundedness scores across languages, suggesting the need for careful study design for groundedness-based research, and the difficulty of grounding verbs in particular.

Chapter 7: Conclusion In this chapter, I summarise the contributions of this thesis, discuss limitations, and outline directions for future work.

1.4 Contributions

Is this needed?

- groundedness
- consistency of lexicality-related distinctions
- connections across levels of linguistic structure
- new computational methods and approaches to typological research

Chapter 2

Background

In this chapter, I argue for a new perspective on comparative concepts in linguistic typology, which grounds operationalizations of complex hybrid concepts in empirical measures of underlying linguistic dimensions. I present the major goals of typology and the challenges of defining comparative concepts for typology, and review existing approaches to these challenges and how they compare to my proposed approach. I then review the development of deep learning models of language, describing how they provide new avenues for defining empirical measures of semantic dimensions of language, and the richness of the semantic and conceptual information they acquire. I then provide a review of the study and application of comparative concepts in *computational* typology, highlighting how building rich empirical models of underlying semantic and perceptual spaces have been key to successful computational typological research, and the parallels between these approaches and my proposed approach. I also highlight the shortcomings of current discrete approaches to semantics in computational typology. Together, this motivates the empirical grounding approach to comparative concepts and linguistic categories that I take in this thesis.

Finally, I provide a high level overview of the lexicality spectrum, defining

formal and functional dimensions of lexicality, and describing their interrelationships. This sets the stage for the remainder of the thesis, which focuses on defining empirical measures of these dimensions, leveraging deep learning models of language, and investigating how these dimensions relate to existing lexicality-related distinctions in multilingual databases.

2.1 Typology and Comparative Concepts

Linguistic typology is the study of variation across the world’s languages. Typologists perform cross-linguistic comparisons with the aim of making generalizations about this variation. Such generalizations may consist of identifying and classifying languages into a small set of types (typological classification) or identifying cross-linguistically consistent patterns in variation. By studying this variation, typologists aim to identify the limits on and universals of human languages, and, often, to identify simple, language-neutral explanations of these limits.

To make cross-linguistic comparisons and identify cross-linguistic variation, typological research has explicitly or implicitly had to identify a frame of *alignment* between languages—typically taking the form of shared concepts identified across languages. Take, for example, the study of basic word order typology:

(2.1) *Paul kisses Peter.*
 SUBJ VERB OBJECT

(2.2) *pooru-wa piitaa-wo kisu-shiteiru*
 Paul-TOPIC Peter-ACC kiss-DO-PRES.CONT
 SUBJ OBJECT VERB

English and Japanese, then, vary in their basic word orders. In English, the verb is preceded by the subject and followed by the object (“SVO order”), while in Japanese, the default ordering is subject-object-verb (“SOV”). This comparison, however, relies on the consistent cross-linguistic identification of

the categories of SUBJECT, VERB, and OBJECT. Such concepts over which cross-linguistic comparisons can be made have been termed COMPARATIVE CONCEPTS (Haspelmath, 2010; Croft, 2016b).

Many of these comparative concepts present serious methodological challenges. It is well known that many categories in linguistics are semantically *motivated*, but not semantically *defined*. Take, for example, the category of subject. While subject's across languages are typically the agents of an action, in the specifics of individual languages, there is additional complexity. While English has "SVO" order, in passive constructions, it is the patient and not the agent which appears in this initial position:

(2.3) *Paul is kissed by Peter.*
 SUBJ VERB OBJECT

Defining the subject in terms of the obvious distributional commonality between *Paul* in (2.1) and (2.3) would make the claim "English is an SVO language" circular. While getting around this particular problem is relatively straightforward (through the deployment of additional constructional tests), this type of issue is pervasive in typological analysis, and careless or inconsistent application of categories cross-linguistically can lead to generalizations or even debates which are vacuous.

For example, frequent debates have occurred over whether a particular language has the category ADJECTIVE: a syntactic category covering property words, distinct from nouns and verbs. The typical structure of such debates involves identifying the behaviour of words which denote properties in a particular language in various constructions, and comparing their behaviour to members of other classes. For example, in Korean, both adjectives and verbs (but not nouns) inflect for tense, leading some to argue that Korean lacks a class of adjectives, and to claim that in Korean, adjectives are a type of stative verb. On the other hand, some have argued that because adjectives in Korean are somewhat

restricted in terms of the tense–aspect–mood constructions they can appear with, they are better analysed as a distinct class. However, cross-linguistically, adjectives rarely inflect for tense or aspect, so this distinction is being made on a very language-particular basis.

Croft (2001) calls this type of syntactic argumentation **METHODOLOGICAL OPPORTUNISM**: the application of arbitrary language-particular criteria to identify distinctions between supposedly universal categories. This approach cannot lead to consistent generalizations across languages. If we consider a generalization like “adjectives do not inflect for tense”, then Korean is a counterexample if adjectives are not a type of verb. If they are a type of verb, then Korean is a counterexample to the generalization “adjectives require some kind of copula-like element in prediction”. To understand what the actual generalizations in typology are and whether a particular language is or is not a counterexample, we need to base these comparisons on cross-linguistically consistent criteria.

What the best comparative concepts are for a given problem is an empirical question, based on their predictive power in terms of generalizations about language variation. All the comparative concepts I have discussed so far are **HYBRID CONCEPTS** (Croft, 2016b): they combine aspects of formal distribution with semantics. As an alternative, we might consider **FUNCTIONAL**¹ comparative

¹By simultaneously addressing both issues of comparative concepts and the lexical–functional distinction in this thesis, I am trapped into a very confusing overload of the term *functional*; it has two, almost diametrically opposed meanings in the literature. As discussed in Chapter 1 the context of the lexical–functional distinction, *functional* refers to a pole of a continuum of linguistic behaviour, where “functional” items/elements/units are those which serve primarily to organize and clarify relationships between other elements. These elements are often described as “grammatical” or “lacking meaning”. In the context of comparative concepts and typological theory, however, linguistic *function* refers to the communicative content of linguistic expressions, as contrasted with its *form*: the specific linguistic realization which conveys a function. In this sense, the function of an adjective is basically its intension: the property it denotes, while the function of a verbal tense marker is the temporal and aspectual information it conveys about the event describe by the verb. Thus, a “functional comparative concept” in this sense is one which is defined in terms of the communicative content of the linguistic expression, rather than its formal distribution. Here, the terms “function” and “functional” are preferred to “semantics” and “semantic” because the later terms are often taken to refer only to truth-conditional content, while function is inclusive information structure.

Of course, these senses are related, but in a quasi-antonymic manner: more functional elements have more abstract, relational, and language-internal functions. In this thesis, when I use the

concepts: concepts that invoke only the communicative function of the linguistic expressions studied, regardless of their formal distribution. Returning to the adjective example, this would be making comparisons of all expressions of property meanings across languages, as suggested by Haspelmath (2012). However, as noted by Croft (2016b), such broad functional comparative concepts may fail to capture important cross-linguistic distinctions that are relevant for typological generalizations. Different types of properties may have different distributions within a single language: the expression of colour properties may pattern differently from emotions, for example. Thus, increasingly fine-grained functional comparative concepts are often necessary to capture cross-linguistic variation. Even with purported “functional” comparative concepts, something of the formal tends to creep in. Ultimately, after all, the typological generalizations are about the behaviour of linguistic expressions, which are formal entities. There are therefore two major dimensions of challenge in identifying useful, valid comparative concepts for typology: selecting the right functions, and formalizing categories of linguistic expressions which align with these functions. In the next section, I review different meta-approaches to these challenges, and describe the general way I tackle these issues in the thesis, which diverges from prior work in important ways.

2.1.1 Defining Hybrid Comparative Concepts

Constructions and strategies ? provides a useful terminological and conceptual framework for describing how typological research can and has often implicitly combined form and function into useful cross-linguistic generalizations, which I will use to frame the discussion in this section. We will follow ? to define a FUNCTIONAL CONSTRUCTION (p. 17):

word “functional” to describe words, morphemes, items, elements, or units, I mean it in the lexical–functional sense. Other uses of the words “function” or “functional” generally refer to the communicative content sense (function as opposed to form), unless otherwise specified.

any pairing of form and function in a language (or any language) used to express a particular combination of semantic content and information packaging²

This type of construction is defined only by *what* it expresses. Croft contrasts this with the narrower STRATEGY ?, p. 19:

a construction in a language (or any language), used to express a particular combination of semantic structure and information packaging function (the *what*), that is further distinguished by certain characteristics of grammatical form that can be defined in a crosslinguistically consistent fashion (the *how*).

The separation of strategies from functional constructions is a useful conceptual tool for approaching conflicting cross-linguistic formal data about the expression of a particular function. For example, one argument that Korean Adjectives are a type of verb is that Korean Nouns require a copula to be predicated, while Adjectives do not. Rather than saying “Korean lacks adjectives”, we can say that Korean uses different strategies for predicating nouns and adjectives, with nouns using a copula strategy, and verbs using a zero-copula strategy. This method of description foregrounds the actual distributional data that needs to be explained. Indeed, empirically there are many interesting questions about the cross-linguistic distribution and co-occurrence of different strategies for particular linguistic functions.

Semantic maps An extremely influential method for relating form to function in linguistic typology is the use of SEMANTIC MAPS (Haspelmath, 2003; Croft, 2002a, pp. 133–139).³ As I discussed, a major problem with replacing traditional comparative concepts like “adjective” with broad functional comparative concepts like “properties” is that languages may have different formal behaviour

²Here, “information packaging” refers to the discourse organization of semantic content within an utterance. It is not of central importance to the present discussion.

³These sources summarize the emerging literature around semantic maps and standardize terminology; the method was developed over gradually over a few decades by a number of linguists, as described in the referenced passages.

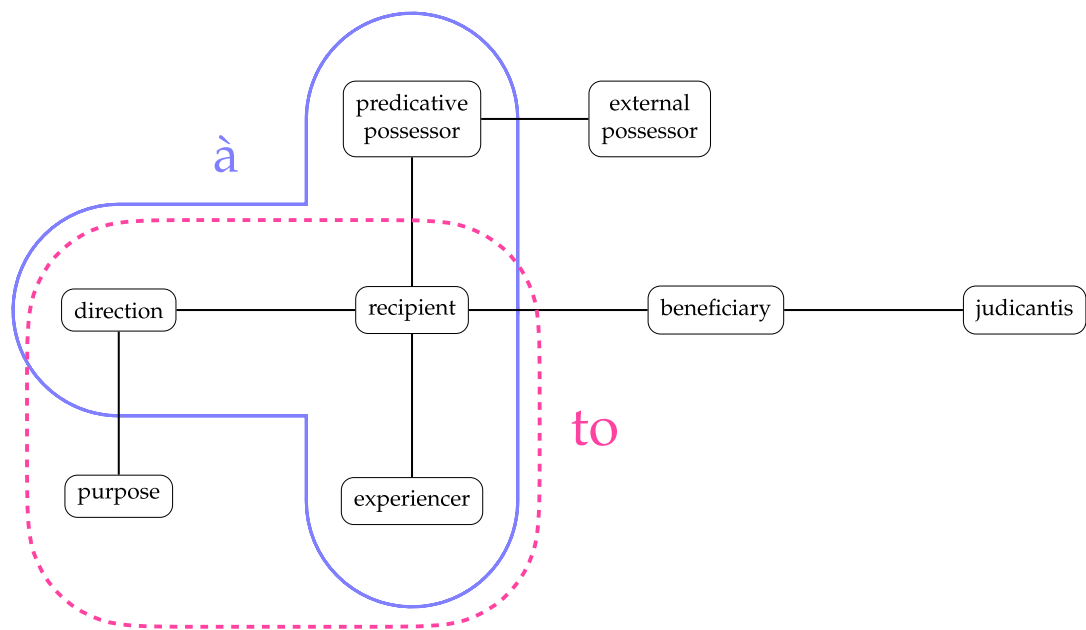


Figure 2.1: An example semantic map for the dative domain, adapted from Haspelmath (2003). Nodes represent different functions which “dative-like” elements can express. The boundaries for English *to* and French *à* are shown in pink and blue, respectively. Both terms cover contiguous regions of the map, satisfying the Semantic Map Connectivity Hypothesis.

for different properties. As such, we may need to define more fine-grained functional comparative concepts. But what if we are still interested in a broader function like “properties”?

Semantic maps offer a solution by decomposing broad functional domains into a network of finer-grained functions, which can then be related to one another based on their observed co-expression⁴ across languages. Figure 2.1 shows an example of a semantic map. To construct such a map, the functional categories are first selected on the basis of whether at least one pair of languages differ in their expression of a function. English uses *to* both for direction (“I’m

⁴That is, using the same form or construction in a language. English co-expresses singular and plural second person as *you*.

going to the store”) and purpose (“I’m leaving early to be on time”), while French uses *à* for direction but *pour* for purpose. Thus, “purpose” must be a distinct functional category. The functional categories are then organized into a graph structure based on co-expression, with the aim to make language-particular strategies correspond to connected subgraphs of the semantic map. This desiderata has been termed the SEMANTIC MAP CONNECTIVITY HYPOTHEIS (Croft, 2001, p. 96), and has also been claimed as a universal property of human language. Designing maps to satisfy this property has the following corollary for the resulting map: if function A and function B are co-expressed in a language, and there is no path between A and B on the map that does not pass through function C, then C must also be co-expressed with A and B in that language. Croft (2001) calls this resulting graph the *conceptual space*, and this is the resulting language universal that is claimed by a particular semantic map analysis. On top of the conceptual space, we draw SEMANTIC MAPS for particular languages, which show the subgraphs of the conceptual space that are co-expressed in that language. The conceptual spaces, then, represent cross-linguistic *constraints* on the application of strategies to express particular functions cross-linguistically, thereby providing a comparative bridge between form and function. Semantic maps can be very useful for expressing the generalizations about problematic, broad hybrid comparative concepts like “dative” or “adjective.”

The semantic map method and the semantic map connectivity hypothesis underlying it have been extremely fruitful in identifying cross-linguistic co-expression patterns and universals. While conceptual spaces are not based on semantic similarity *per se*, rather facts of co-expression, in many domains where the semantic map method has been applied, the resulting conceptual spaces align closely with semantic similarity. In Section 2.3.4, I will describe how computational methods have built on the theory and practice of semantic maps, as well as some of the limitations of existing approaches.

Retro-definitions Haspelmath (2021) proposes a radical approach to hybrid comparative concepts, which he terms **RETRO-DEFINITION**. In this approach, common but potentially problematic comparative concepts (“traditional comparative concepts”) are maintained as terms, but re-defined in a way that is cross-linguistically straightforward to operationalize and closely matches the traditional term. In this way, it represents a radical acceptance that comparative concepts need not relate to any “true” categories of language. As an example, Haspelmath proposes retro-defining adjectives as “property roots”, regardless of their syntactic behaviour in a given language. Even more radically, he proposes retro-defining “inflection” as morphemes that express a fixed set of meanings cross-linguistically, and “derivation” as any kind of word-formation process that expresses any other type of meaning (Haspelmath, 2024). This has the effect of allowing the precise usage of these terms in cross-linguistic comparison, but whether these definitions provide the most *useful generalizations* about language is an open question.

Prototype theory and fuzzy categories Another conceptual approach to dealing with the challenges of defining comparative concepts is to embrace the idea that categories are inherently fuzzy and gradient, and organized around a central *prototype*. This viewpoint was popularized in cognitive science by a series of seminal works by Eleanor Rosch, which demonstrated clear effects that people both agree which members of categories like “vegetable” or “furniture” are most prototypical, and that prototypicality influences processing (?). This finding was influential on early work in the development of cognitive linguistics, with theorists like Ronald Langacker and George Lakoff arguing that linguistic categories also have a prototype structure ??, and has been proposed as a solution to conflicting cross-linguistic data about categories: unusual distributional behaviour is associated with less prototypical members of a category.

However, the approach has come under fire for failing to account for apparent category boundaries or being unfalsifiable when applied to distributional data (Haspelmath, 2024).

In many instances, prototype models are like a crude version of a semantic map model, because the semantic map connectivity hypothesis has similar implications: the longer the path between functions in conceptual space, the less likely they are to be co-expressed. However, rather than providing a fine-grained map of functions, prototype models focus on identifying central features, and suppose that less prototypical members should be less likely to be co-expressed.

My approach: grounding comparative concepts The main approach in this thesis does not fall neatly into any of the above widely-discussed approaches to hybrid comparative concepts, but takes a heterogeneous set of inspirations from them. The approach of the thesis is to define empirical measures which capture continuous dimensions of formal and functional variations, and relate them to *existing* category operationalizations. While this approach will be described in more detail later in the thesis, especially in Chapter 4, I will here briefly review how it relates to these existing approaches. Similar to the semantic map approach, I aim to take a complex, hybrid concept and decompose it into finer-grained dimensions. Because my measures are continuous, it shares with many versions of the prototype approach the idea that category membership is gradient. My approach takes from retro-definitions the idea that comparative concepts come from a linguistic tradition which needs to be questioned. However, rather than re-defining existing terms, I take an existing operationalization as a starting point, and investigate how well my empirical measures align with these operationalizations. In the thesis, I focus on standard, theory-neutral operationalizations from databases like UniMorph and Universal Dependencies. Because the relationship between my empirical measures and the operationaliz-

ations is the object of study, the claim does not rest on the a-priori validity of the databases—we aim to study the relationship of these comparative concepts *as they are used*. Future work can and should see if different operationalizations align better with the empirical measures I define. This empirical grounding approach also takes inspiration from literature in computational typology, which I will review in Section 2.3, but to my knowledge this literature has not been directly invoked in the context of comparative concepts.

Debates around comparative concepts often focus on problematizing existing attempts (Croft, 2016b, p. 379). With this approach, I aim to reverse the discussion, by quantifying the consistency of existing operationalizations with respect to my empirical measures. In the context of this thesis’s focus on *lexicality*, this means relating distinctions among word classes or between inflection and derivation to highly *abstract* empirical measures; in contrast to the semantic map approach or many of Haspelmath’s retro-defintions, which focus on the *functions themselves*. In this way I aim to provide a new perspective on whether such distinctions really relate to the abstract properties which are often invoked to explain them, but have been difficult to measure directly. To do so, I rely heavily on emergent computational techniques for learning linguistic representations, which I will now describe in the next section.

2.2 Finding Meaning in Computational Models

Over the past two decades, deep learning models have revolutionized the field of natural language processing. These models come out of a history of brain-inspired cognitive modelling called *connectionism* (?), which posited that important aspects of human cognition were best understood as the emergent behaviour of large parallel and distributed networks of simple processing units. In natural language processing, deep learning models have demonstrated the

ability to perform linguistic tasks at a level that would have once been thought to require human-level linguistic competence, like translating sentences or summarizing documents (?). While these models do not in-and-of-themselves provide a theory of human language processing, they provide a useful test bed for gradient and usage based theories of language (?). In this section, I discuss several types of models, and the evidence that they acquire rich semantic and conceptual representations.

2.2.1 Distributional Embeddings

Describe how they learn meaning from co-occurrence, in line with the distributional hypothesis. Discuss evidence of semantics, measurement (similarity, directions). Bilingual lexicon induction as evidence for a shared conceptual space

2.2.2 Contextual Embeddings and Language Models

Discuss the type-level limitations of distributional embeddings, and how contextual embeddings address these. Evidence of semantics, average embeddings vs. contextual. Language models, specialization for prediction, semantic evidence (lexical, structural).

2.2.3 Multimodal Models

Disentangling form and function. How this relates to typology. How they work, grafting a vision model onto a language model. How the separate modality provides an avenue for language-agnostic function. universal representation hypothesis.

2.3 Approaches to Comparative Concepts in Computational Typology

This section reviews the roles that comparative concepts have played in computational approaches to linguistic typology. In this section, I aim to highlight how technological advances can enable new types of comparative concepts, and how the choice of comparative concepts can shape the types of questions that can be asked and generalizations that can be formed. In particular, I will argue that computational modelling thrives when fine-grained distance metrics are defined, be this in a high-dimensional discrete space or especially a continuous space.

2.3.1 Comparative concepts in multilingual databases

The predominant approach to comparative concepts in computational typology has simply been to follow the lead of whatever annotation scheme is used in the typological or cross-linguistic database used in the study. For example, influential works like ?, which demonstrated minimization of dependency length as a universal pressure on word order, used the data of Universal Dependencies (?), as-is. Similarly, cross-linguistic studies of morphological complexity and inflectional paradigms have relied on a combination of feature values from databases like the World Atlas of Language Structures (WALS) (?), or the encodings of grammatical features in UniMorph (Batsuren et al., 2022).

Databases like UniMorph and Universal Dependencies attempt to use a single cross-linguistic annotation scheme for their comparative concepts, but these schemes usually much more closely resemble the hybrid categories of language-particular analyses. For example, UniMorph’s feature set includes values like `TENSE=PAST`, with no information about what constructions that form is used in. Chapter ?? includes a detailed discussion of some of the limitations

of UniMorph annotations, which are largely based on language-specific grammatical traditions rather than typological best practices. Similarly, Universal Dependencies part-of-speech tags are cross-linguistically “universal”, but they are deployed with the “methodological opportunism” described by ?, where categories are defined in a language-particular way. Very recently, Universal Dependencies has begun a process of aligning their representational scheme with more fine-grained and cross-linguistically valid comparative concepts (??), but this process is still ongoing.

That the annotation of these databases is *not* consistently based on cross-linguistically valid universal comparative concepts actually provides an interesting opportunity, however, and one which this thesis exploits. A major aim of this thesis is to define computational comparative concepts of dimensions hypothesized to underlie grammatical category distinctions cross-linguistically, and then investigate how these dimensions relate to the categories used in these databases. In this way, I can investigate the extent to which these databases distinctions, while flawed, nevertheless align with the underlying dimensions of meaning that motivate these distinctions.

2.3.2 Phonological typology

Phonology offers the earliest and clearest example of how empirical, continuous measures can advance typology. It has long been recognized that language-specific categories like phonemes are ill-suited for cross-linguistic comparison (“phonemes are not fruitful universals”; ?) because they are defined by language-internal contrasts. To address this, ? introduced a *feature-based* account of phonetic universals, later refined by ? through the notion of *natural classes*.

Vowels have long been central to typological inquiry. ? proposed early implicational universals about vowel systems within his proto-featural framework. Yet within this featural approaches, universals of vowel systems were

understood as complex, varying substantially on the number of vocalic contrasts within a language. The true generalizations, which turn out to be extremely simple when properly understood in a continuous space, required substantial developments in the acoustic theory of vowels. The decisive shift came with the formant theory (??), which linked vowel quality to acoustic resonances (F1–F3). The development of better technologies for measuring and recording formants through the twentieth century ultimately provided a real-valued acoustic representation (?) that allowed vowels to be compared across languages in a shared empirical space, enabling the development of theories that made quantitative, testable predictions. The quantal theory (??) proposed that languages prefer perceptually stable regions of this space, while the dispersion theory (?) modeled vowel inventories as systems maximizing perceptual distance. By simulating optimal vowel systems and comparing them to attested inventories, these models offered a precise computational account of typological tendencies. Subsequent dispersion–focalization models (??) and probabilistic analyses of entire vowel-system distributions (?) further refined these predictions, revealing the relative influence of competing pressures such as distinctiveness and perceptual stability.

While vowel typology concerns formal acoustic dimensions, its trajectory exemplifies a broader lesson: defining an empirical, continuous underlying space can transform typological theory. Just as formant space enabled precise and falsifiable generalizations about vowel systems, deep learning models may provide an analogous empirical grounding for meaning. Linking model-derived semantic spaces to human cognition could allow typology of linguistic function to achieve the same level of quantitative precision. In the next section, I turn to evidence from semantic category systems supporting this view.

2.3.3 Semantic category systems

The domain most closely paralleling vowel typology in its treatment of function is the study of semantic category systems—cross-linguistic analyses of how languages partition a shared underlying semantic space. As with vowels, researchers model these systems as optimizing trade-offs among universal pressures such as simplicity, communicative efficiency, and learnability. Once an underlying space is defined, these pressures can be formalized, simulated, and quantitatively tested against attested systems.

The case of color terms provides the clearest illustration. In their seminal study, ? identified robust implicational hierarchies—two-term systems distinguishing light from dark, three-term systems adding red. Berlin and Kay couched these generalizations in terms of an implicational hierarchy of color terms; however, their methodology could offer little insight into *why* these hierarchies exist. The breakthrough came with the development of a precise underlying perceptual space for color. Studying the relationship between the color space in terms of frequency and perceptual distance, Later work showed that a continuous perceptual space was key. Building on the CIEL*a*b* color space (?), which aligns physical and perceptual properties of color, ? modeled how systems maximize within-category similarity and between-category distinctiveness. These simulations closely predicted attested color inventories and revealed that real systems are significantly more optimal than chance. This initial effort has expanded into a rich literature, modelling trade-offs among different pressures such as communicative need, perceptual structure, and learnability. While this remains a more rapidly evolving area of research than vowel system typology, studies continue to refine hypotheses and distinguish between pressures with increasingly fine-grained predictions about color systems.

Comparable approaches have since been applied to other semantic domains—kinship, number, quantifiers, modals, and pronouns—where discrete

conceptual structure allows for tractable mapping. More challenging are domains with continuous meaning spaces, such as spatial terms (?) and tense-aspect marking (?), which have required simplifying the space into coarse discrete categories or low-dimensional projections (e.g. multidimensional scaling of usage data; ?). Therefore, work in these frameworks has primarily focused on domains where an underlying conceptual space can be more straightforwardly defined.

Overall, both the study of vowel systems and semantic category systems show that by creating an operationalization of an underlying space, we can test and discover parsimonious and predictive theories of the fundamental data of linguistic typology. In each case, substantial typological progress was able to be made without access to the “true” space—we are still learning about the perceptual dimensions of vowels (), and CIEL*a*b* has known shortcomings ()—but the development of some empirical model of the underlying space that could be coded across languages was critical for this progress. However, the limited scope of functions studied in these types of computational frameworks points to the challenge of defining an underlying space for more complex semantics. In this thesis, I argue that recent advances in deep learning models of language are likely an early step in this direction, analogous to the early stages of development of the formant theory of vowels.

2.3.4 Multidimensional scaling

The multidimensional scaling (MDS) approach to semantic maps proposed by ? represents the most well-developed technique for modelling a continuous functional space for typological comparison. This approach was developed to address the challenge of translating large typological datasets into a traditional

semantic map following the methodology of ?,⁵ and to provide a stronger mathematical basis for the semantic map theory. These methods take as input a high-dimensional discrete matrix of linguistic data, and produce a low-dimensional continuous representation of the data, using either optimal unfolding, or in some studies, matrix decomposition techniques. The resulting map provides an approximation where Euclidean distance approximates the frequency with which two functions are expressed with the same form. In this way, an underlying semantic/conceptual space is being inferred from typological data about form—representing a move away from the discrete identification of functional comparative concepts to a richer emergent empirical representation of the complexities of meaning, towards the desiderata of this thesis. Studies vary primarily in the way they construct the input dissimilarity matrix, and thereby in how much they allow for an emergent representation of function. ? provide a recent overview of the different methods used in the literature, which I briefly summarize here to illustrate how differing approaches rely on different comparative concepts. They provide a three-way typology of approaches to constructing the input for an MDS analysis.

First, there is the classical input representation, which ? used to recreate ?'s analysis of indefinites. Other examples of this type of map include... This type of map takes a set of N linguistic forms \mathcal{F} , and a set of K underlying functions \mathcal{M} . The binary input matrix $\mathbf{I} \in \{Y, N\}^{K \times N}$ is constructed such that

$$\mathbf{I}_{i,j} = \begin{cases} Y & \text{if form } f_i \text{ from language } l \text{ conveys (or is used to express) function } m_j, \\ N & \text{otherwise.} \end{cases}$$

to which the optimal unfolding technique is applied. In this type of analysis, the functions are entirely manually posited by the typologist and abstracted away from the constructions on which the functional claim is based. As a result,

⁵While an algorithm for producing graph-based maps from large-scale data was later introduced by ?, the MDS techniques retain a number of advantages.

the MDS analysis cannot capture fine-grained variations around the prototypes of a given abstract function, but can capture differences in the closeness of two functions in a more fine-grained way than the classical approach to semantic maps—frequency of form-function co-occurrences is modelled. Nevertheless, in terms of comparative concepts, the functions here retain the traditional approach to function in typology, with its known shortcomings.

? also introduce a second method for producing an MDS map, which allows more fine-grained study of the prototype structure of functions. This second map relies on ?'s tense-aspect data. Here, rather than manually determining in a binary manner whether a particular linguistic form can or cannot encode a particular function, a range of specific constructions are included in the analysis. In ?, informants across languages translated sentences in a specific temporal and observational context (e.g., you saw someone writing a letter yesterday). These data were assigned to tense-aspect prototypes, so the input matrix \mathbf{I} now has K sentential contexts $c_j \in C$, belonging to a smaller number of abstract function prototypes, and takes the following form:

$$\mathbf{I}_{i,j} = \begin{cases} \text{Y} & \text{if form } f_i \text{ was used for sentential context } c_j \text{ in language } l, \\ \text{N} & \text{otherwise.} \end{cases}$$

The lessened reliance on manually posited functions allows for a the prototype structure to emerge from the data. However, the contexts are still manually selected by the typologist, and the semantic information still only comes from co-occurrence with forms in the sample itself—so the study necessarily cannot represent the full complexity of the studied forms across the vast space of possible meanings, nor can it fully capture frequency effects. This type of study has also been applied to other aspectual constructions (?), and to verb-specific semantic roles (?).

Finally, the third major type of MDS analysis relies on a fully bottom-up approach to function, using parallel corpora rather than reference grammars

or elicited survey data. In this type of study, relevant parallel clauses for a particular phenomenon are identified in a parallel corpus, and the input matrix \mathbf{I} is constructed such that each row contains the construction used in that clause in each language studied, producing a K -tuple where K is the number of languages. To compute distance in this type of study, the Hamming distance between the tuple for two clauses is computed, yielding a similarity matrix based on the number of languages that use the same construction in both clauses. This type of analysis removes the manual positing of functions and contexts entirely, proceeding bottom-up, and is thus the most in the spirit of our present inquiry. However, the approach still only captures similarity based on translations in corpora. The fact that contemporary deep learning models capture extremely fine-grained semantic distinctions is not leveraged in this approach, and so the semantic space captured is only as good as the evidence directly given by the co-occurrence of translations, rather than language-internal evidence about meaning.⁶

Overall, the MDS approach allows us to both study the rich gradient structure underlying linguistic function, and decrease the dependence on manually posited functions. However, the state of the art still relies entirely on parallel co-occurrence data. In the next section, I will discuss the small literature that leverages recent advances in deep learning to provide a rich representation of function in typology.

2.3.5 Deep learning models of comparative concepts

Despite the rich body of evidence that deep learning models of language capture fine-grained semantic distinctions, there has been relatively little work leveraging these models to provide empirically grounded comparative concepts for typology. Recently, ? used multilingual BERT (?) and Aya (?) to study animacy

⁶A wide range of domains and phenomena have been studied with this approach:

cross-linguistically. Specifically, they identify which syntactic roles and clausal positions are most associated with animacy of the referent. However, the role of the models used here is not truly gradient, nor is the function emergent—the models are used to produce a 3-way classification (human, animate, and inanimate) based on an annotated corpus, and the analysis is conducted over these discrete categories. While the rich representations of the models are critical for creating an accurate classifier, the comparative concepts are still discrete and manually posited.

? study grammatical subjecthood with a less discrete approach. Specifically, they train a multi-layer perceptron classifier on multilingual BERT representations to distinguish between the embeddings of transitive subjects and objects, then examine the classifier’s categorization of intransitive subjects, finding that intransitive subjects are categorized as more subject-like than object-like, and that classifiers transfer across languages, including languages with different morphosyntactic alignment (e.g. ergative-absolutive vs. nominative-accusative languages). However, they found that animate non-subjects and passive subjects were more likely to be classified as subjects and objects respectively, indicating a semantic dimension to this cross-linguistically robust representation of subjecthood.

Another technique for obtaining a gradient representation of a semantic dimension is to use *semantic projection* (?), which has been shown to capture human judgements about object features. This technique uses exemplars at extremes of a semantic dimension (e.g. “huge” and “tiny”), using a deep learning model to embed them in a Euclidean space. All possible embedding pairs across the two sets of exemplars are subtracted from each other, and these difference vectors are averaged to produce a single vector representing the semantic dimension. This vector can then be used to project other words onto this dimension by taking the dot product of their embedding with the dimension

vector. ? uses this technique to study models' representations of animacy. The authors claim that their results show that models represent animals as more animate than humans, in line with psychological findings in humans (). They suggest that this indicates inductive biases in humans that shape grammatical animacy by focusing on certain constructions. However, their operationalization of animacy is questionable, as the exemplars they use to define high animacy are exclusively non-human animals. Nevertheless, the techniques here show how deep learning models can be used to provide a gradient representation of a semantic dimension which can be used to study cross-linguistic patterns in form-function mappings.

Altogether, the results in this nascent literature are promising, but there are still many dimensions of deep learning representations that have not been explored. In this thesis, I will focus on how deep learning models help provide new and better models of lexicality, which has so far not been directly addressed in any of this literature.

2.4 Formal and functional dimensions of lexicality

...we may be quite sure of the analysis of the words in a sentence,
and yet not succeed in acquiring that inner “feel” of its structure
that enables to tell infallibly what is “material content” and what is
“relation”

— Edward Sapir (?)

Some units of language are more meaningful than others. This basic insight is almost as old as the study of language itself. In the Greek tradition, Aristotle distinguished *phōnē* (sign-bearing sounds) from *phōnē ásemos* (non-sign-bearing sounds), such as the class of *árthron* which includes prepositions and preverbs (?). This distinction was not limited to the proto-linguistics of Indo-European languages: in the 12th century the *Wén zé* (文則) of Chen Kui (陳葵) catalogued *zhùchí* (助詞) (lit. “helping words”)—corresponding to what we would

today call function words. Across the world's languages, we see asymmetries between elements that express content and those that express grammatical function. It is little wonder then that the distinction between contentful and functional elements continues to have relevance across linguistic theories and domains. Yet boundary cases abound and the nature of the distinction has made it challenging to formalize. In Chapter 1, I sketched the idea of a **LEXICALITY SPECTRUM** as a general way of conceptualizing the correlation between formal expression and (degree of) semantic content. In this section, I describe the formal and functional dimensions of this spectrum in more detail.

2.4.1 The formal dimension

As I argued in Chapter 1, the lexicality spectrum is operant at multiple levels of formal linguistic structure, with different names being used for similar distinctions at different levels. But the distinctions between “words”, “clitics”, “morphemes”, and “affixes” are all theoretically tenuous at best (Zwicky, 1994; Bruening, 2018). As a theory of these terms themselves is beyond the scope of this thesis, I will here focus on the general formal trends that underlie these distinctions—the formal dimension of the lexicality spectrum.

The basic idea underlying all these terms is this: some linguistic units have “bigger” forms than others. Of course, this is implied by the compositional nature of linguistic structure: a phrase may be composed of several words, a word may be composed of several morphemes, and each morpheme can contain a variable number of phonemes. Yet even comparing morphemes to morphemes, some are formally bigger than others. Here are some ways in which this can manifest:

Boundness One aspect of formal size is the notion of **BOUNDNESS**. While Haspelmath has argued for a sharp cross-linguistic definition of boundness ?? as

“unable to occur in isolation”, I share ?’s scepticism of the utility of this as a cross-linguistic criteria and share his feeling that this is better understood as a gradient notion. There are many languages where no morpheme can occur in isolation—surely our comparative concepts should apply to them! Further, the notion of “isolation” is itself problematic, as language always occurs in a discursive context. Here, I sketch boundness as a continuum of cluster properties. In Chapter 3, I operationalize some relevant aspects of this continuum, but here I will simply focus on what the formal trends *are*.

At one end of the spectrum of boundness are free morphemes. In many languages, these morphemes can form whole utterances by themselves in the right context (Consider “Cat.” as a response to “What is your favourite animal?”). In some languages, even the freest morphemes may not be able to stand alone, requiring some obligatory bound marking (e.g. case or tense marking), but the free morpheme behaves in some way like the “root” of the word. This often takes the form of the free morpheme occurring at the periphery of the word (usually the beginning).⁷ Further, they typically occur immediately coincident to that host morpheme. If morphemes occur between a bound morpheme and its root, those morphemes are typically intermediate in terms of these formal properties—that is, the most bound morphemes occur furthest from the root morpheme.

Boundness is a similar notion to VALENCY, the notion that certain linguistic units require a certain number of arguments. For example, nouns typically have a valency of zero. Most verbs, on the other hand, have a valency of one or more, requiring a subject and one or more objects. Syntactic valency is different from semantically valency. For example the verb *to rain* arguably requires no semantic arguments in English, but it still requires a subject syntactically

⁷In cases where the free morpheme cannot stand alone, a critical aspect of the argument for the *freeness* of this morpheme is typically semantic. However, I am here focusing exclusively on the *formal* properties of boundness.

(*It rains*). Putting things in mathematical terms, we can view such words as functions, which require certain arguments to form a complete expression. Similarly, bound morphemes are often formalized as predicates which take in a root morpheme to produce a combined expression. However, prototypical free but valent morphemes (e.g. verbs) are distinguished from bound morphemes by their degree of syntagmatic integration with their arguments. Either the arguments are themselves free morphemes, able to move around depending on the construction or take their own bound morphemes, or else they are expressed through bound morphemes on the verb itself.

Allomorphy More bound forms are also phonetically more variable. They may be subject to special phonological processes that do not apply to other morphemes in the language (such as the English plural -s being realized variably as [s], [z], or [ɪz] depending on the phonological context). The more of these processes apply, the more phonologically bound the morpheme is. ? argue that lexically- or morphologically-conditioned variants of morphemes (*allomorphs*) are a formal sign of greater bondedness.

Length Perhaps the most obvious dimension of formal size can be seen in the number of phonemes in a morpheme—which can vary dramatically. The length parameter, at the short end, is intimately tied up with the other dimensions of formal size. Allomorphy may reduce the number of shared phonemes between morpheme variants. Morphemes can also get a length shorter than one phoneme in some instances. *Portmanteau* morphemes share multiple (unrelated) features in a single marker, meaning that the phonological material dedicated to any one of them can be the equivalent of less than a single segment. Tightly bound morphemes can become shorter than a segment by becoming suprasegmental or process morphemes, e.g. by changing tone or root morpheme vowel quality (*Ablaut*). We can therefore think of all dimensions of formal size outlined here

as related to the concept of length.

2.4.2 The functional-semantic dimension

Information and Frequency At the semantic core of the lexicality spectrum is the notion of CONTENTFULNESS.⁸ The basic intuition is that some linguistic items contribute more to the overall meaning of an utterance than others. This is a notoriously difficult notion to pin down, as it relates to the deep question of what *meaning* and *content* are in the first place. INFORMATION THEORY both provides a mathematical formalization of this notion and a demonstration of why separating content from form is difficult. Shannon (1948) introduced the notion of *entropy* as a measure of information in bits. The entropy of a random variable X with possible values $\{x_1, x_2, \dots, x_n\}$ and probability mass function $P(X)$ is defined as:

$$H(X) = - \sum_{i=1}^n P(x_i) \log P(x_i).$$

The entropy is related to the notion of “information” in an important respect: it provides a lower bound⁹ on the number of bits needed to encode the value of X . That is, it tells you the most efficient encoding/compression of X must be at least $H(X)$ bits long. Shannon’s entropy and information-theory more broadly have been widely and successfully applied to linguistic theory (). With respect to “words”/“morphemes”, the information-theoretic perspective implies that the information content of a morpheme is its negative log-probability.¹⁰ Under this perspective, more frequent morphemes carry less information. Combining this with the idea from Section 2.4.1 that the formal dimension is all related to length or formal size, we can see the lexicality spectrum as a generalization of ZIPP’s LAW OF ABBREVIATION, one of the most famous and foundational discoveries in all

⁸The terms “semantic weight” and “semantic force” are also common

⁹Technically, an average lower bound over many samples.

¹⁰This follows from the assumption that morphemes are generated independently, which is not true, however, later work has extended the information-theoretic insight to account for dependencies, with a similar overall conclusion (Piantadosi et al., 2011).

of computational linguistics (?), which states that more frequent words are much shorter than less frequent ones. The information-theoretic corollary of this is that more informative morphemes are longer than less informative ones. Frequency certainly plays a central role in the structure of language and the lexicality spectrum. It has been argued to be a driving force in grammaticalization (?), and the so-called ICONICITY OF COMPLEXITY (itself closely related to the notion of a lexicality spectrum) has been argued to be an effect of frequency pressures on efficient communication (?). Nevertheless, debates around the nature of the lexical–functional distinction, inflection–derivation distinction, and related phenomena continue unabated, indicating that frequency and thus standard information-theoretic notions of content are insufficient as a full account of the relationship between meaning and function.

Abstractness and Polysemy Another common way of characterizing the functional dimension of the lexicality spectrum is in terms of ABSTRACTNESS. Functional elements are often described as being more abstract than lexical elements; on the other hand, functional elements are also often more POLYSEMOUS than lexical elements (Haspelmath, 2003). The difference between abstractness and polysemy is not always clear. In Figure 2.1, it is fairly clear that *to* has multiple senses, such as the purpose sense in *I went to the store to buy milk* and the recipient sense in *I gave the book to Mary*. Haspelmath (2003) points out that data is often ambiguous between a *monosemic* position where there is a single vague, abstract meaning that interacts with contexts to serve different functions, and a *polysemic* or even *homonymic* interpretation where there are multiple more specific meanings which share a surface form for motivated or ideosyncratic reasons. Under such an interpretation, the meaning of functional elements is not vague at all, and perhaps not very abstract. Nevertheless, words like *in* which have spawned whole literatures in linguistics and psychology attempting

to characterize their use, and which seem to cover a continuous space of meaning rather than discrete scenarios, are more straightforwardly characterized as abstract than polysemous.¹¹

Another problem for abstractness of meaning is that lexical elements can also be highly abstract (e.g. *idea, angry*) and some traditionally functional elements can be thought of as being fairly concrete (e.g. plural markers). This problem is sometimes waved away by invoking prototypicality, but such an account leaves some serious questions unanswered (Croft, 2002a, p. 225). Croft notes that *concrete* nouns which are frequent do get shorter forms (e.g. *dog, car*), indicating that length is more a function of predictability than abstractness. Nevertheless, while these forms are shorter, they do share many formal properties with prototypical concrete lexical items, such as their lack of boundness. So there is still something *different* about these forms compared to functional items, even if abstractness is not the right way to capture it.

Relationality Perhaps one way out of this conundrum is to focus on semantic RELATIONALITY rather than abstractness *per se*. In the simplest terms, a relational meaning is one that inherently implies the existence of at least one other entity (Croft, p. 67). For example, the concept ROUND is relational, as roundness can only be defined with respect to some entity. On the other hand, CIRCLE is non-relational, despite referring to the same properties. In my view, a key property of relationality of meaning is that when composed with entities, the resulting meaning is less abstract than the relational meaning alone. For example *round* is more abstract than *a round rock*. This can help explain some of the more “concrete” grammatical functions: while plural *forms* are not very abstract (e.g. *cats* is likely to be very concrete), the plurality is itself highly abstract.

There are several barriers to associating relationality with the functional

¹¹I refer here to the spatial sense(s) of *in*.

dimension of the lexicality spectrum. A first objection is that because verbs and adjectives are both traditionally considered lexical and relational, relationality cannot be the defining property of functional items. In Part II of this thesis, I will argue that relationality is a key dimension for shaping the lexicality spectrum, and that this *includes* adjectives and verbs as in some sense closer to functional elements than nouns. I believe the definition of relationality I have given here helps explain this: prototypical *lexical* relational concepts such as GIVING, ROUND, or RED are, I believe, more concrete and more full of intension than functional relational concepts such as PLURALITY or ANIMACY. A key correlate of this objection is confusing relationality with semantic valency. The number of entities required to specify a meaning (or the number of syntactic arguments) is *not* the sense of relationality I am referring to here. Valency manifests iconically in syntax, but relationality as I define it is more closely related to notions like boundness. For example, RED is less relational than PLURAL, despite the fact that both are monovalent, because RED is more conceptualizable on its own.

A second problem for this view is that meaning itself is relational. This view goes under the CONCEPTUAL ROLE THEORY of meaning (Block, 1998) in philosophy of mind, and also undergirds structuralist and distributionalist views of meaning in linguistics (??). Piantadosi and Hill (2022) provide a useful example, pointing out that despite a clear prototypical concrete referent, the concept of POSTAGE STAMP is fundamentally relational, and we can easily imagine *virtual* postage stamps so long as the abstract referent fulfills the role of tracking payment for delivery. Plausibly, this tension could be resolved by countering that the conceptualization of POSTAGE STAMP is not as relational as the inferential or functional extensions it affords—that is, we imagine a concrete, bounded, non-relational object, which has been selected by the complex networks of meaning in our minds to fulfill a relational role. Nevertheless, I think this is a serious challenge to the entity–relation distinction which should be further investigated,

but stands outside the scope of this thesis.

Lastly, a key objection to both relationality and abstractness is that they are hard to specify and often subjectively defined. On this point I agree. For example, the sense that RED is somehow less relational than ROUND, or PLURAL is less relational than GIVING, despite the higher valency of the latter, is intuitive, but difficult to give precise criteria for. This is why a key goal of this thesis is to provide empirical and computational tools for investigating these notions rigorously. While the above discussion is theoretical and subjective, and one can dispute the abstractness I assign to various concepts, the empirical facts in the rest of the thesis remain. The argument here should be seen as a motivation and interpretive lens for the empirical work that follows.

2.4.3 Summary

In this section, I have attempted to sketch the separate formal and functional dimensions of the lexicality spectrum. On the formal side, I have argued that boundness, allomorphy, and length are all correlated dimensions of formal size. On the functional side, I have argued for the importance of relationality alongside information content, and discussed how the former relates to boundness iconically, and the latter relates to formal size via economy principles. Together, this provides a high-level overview of the lexicality spectrum, from nouns at one extreme to inflectional affixes at the other. In the rest of this thesis, I will explore specific aspects of this spectrum and distinctions drawn along it in more detail. Therefore properties specific to inflection and derivation, for example, will be discussed in the relevant chapters.

2.5 Chapter Summary

This chapter has provided a theoretical background for the thesis. In ??, I reviewed the problem of comparative concepts in linguistic typology, discussing approaches like semantic maps and retro-definitions, comparing and contrasting them with my approach of *grounding* problematic distinctions in empirical measures. My approach allows for understanding how consistent existing or proposed operationalizations of comparative concepts are in terms of empirical dimensions. In ??, I provided an overview of a decade of advancements in deep learning models of language, focusing on how these models capture fine-grained and potentially universal semantic distinctions, and how large-scale pretraining is useful for the acquisition of rich linguistic structure. I argued that these models provide a promising avenue for separating meaning from frequency in typology. In ??, I connected the computational literature in typology to the notion of comparative concepts, which have largely been implicit in this literature. From phonological typology and semantic category systems, I argued that computational models have been able to advance typological theory and understanding as technologies develop that allow the empirical grounding of the underlying space, further motivating my approach.

Finally, in Section 2.4, I provided a high-level theoretical overview of the lexicity spectrum, separating the formal and functional dimensions, and arguing for the role of informativity and relationality in shaping the formal dimension. In the rest of this thesis, I will use deep learning models to operationalize and investigate specific lexicity-related phenomena, showing consistent patterns across distinctions that have been difficult to formalize in the past.