# Discovering Latent Concepts in BERT

Fahim Dalvi*  Abdul Rafae Khan*  Firoj Alam   Nadir Durrani   Jia Xu   Hassan Sajjad

{faimaduddin,fialam,ndurrani,hsajjad}@hbku.edu.qa     {akhan4,jxu70}@stevens.edu

## 1. Introduction

- Current research on interpreting Deep NLP models is limited to pre-defined concepts
  - Classical NLP tasks (POS, NER, Chunking etc)
  - Ignores what latent concepts are learned by the model
- We propose a method to analyze latent concepts learned in pre-trained models
  - What are the novel concepts learned by the model?
  - How much do these latent concepts align with pre-defined linguistics concepts?
  - How do concepts evolve across the network layers
- We annotated latent concepts in BERT and provide a multi-facet hierarchical conceptNet dataset (BCN)
  - 174 fine-grained concepts and a total of 1M annotated instances
  - The dataset enables model-centric interpretation
  - The dataset can be used as a new classification dataset for NLP in general

## 2. Methodology

- **Concept**
  - represents a notion and can be viewed as a coherent fragment of knowledge
  - a group of words that are meaningful e.g. *Names of ice-hockey teams, First words of a sentence, Words that begin with "anti"*

- **Methodology**
  - Given a pre-trained model and a corpus of sentences
  - Extract contextualized representations of words
  - Group words into clusters using hierarchical clustering
  - Manually annotate each cluster into fine-grained categories
  - Analyze the cluster by aligning them with the pre-defined linguistic concepts

## 3. Annotation Task



- We annotated 279 clusters
  - 243 (87.1%) meaningful and 36 (12.9%) non-meaningful clusters (Q1)
  - 142 (75.9%) can be combined with the sibling to form a bigger meaningful cluster (Q2)

- BCD dataset consists of 174 unique concepts
  - 11 lexical labels, 10 morphological labels, 152 semantic labels, and 1 syntactic label

## 5. Annotated Dataset

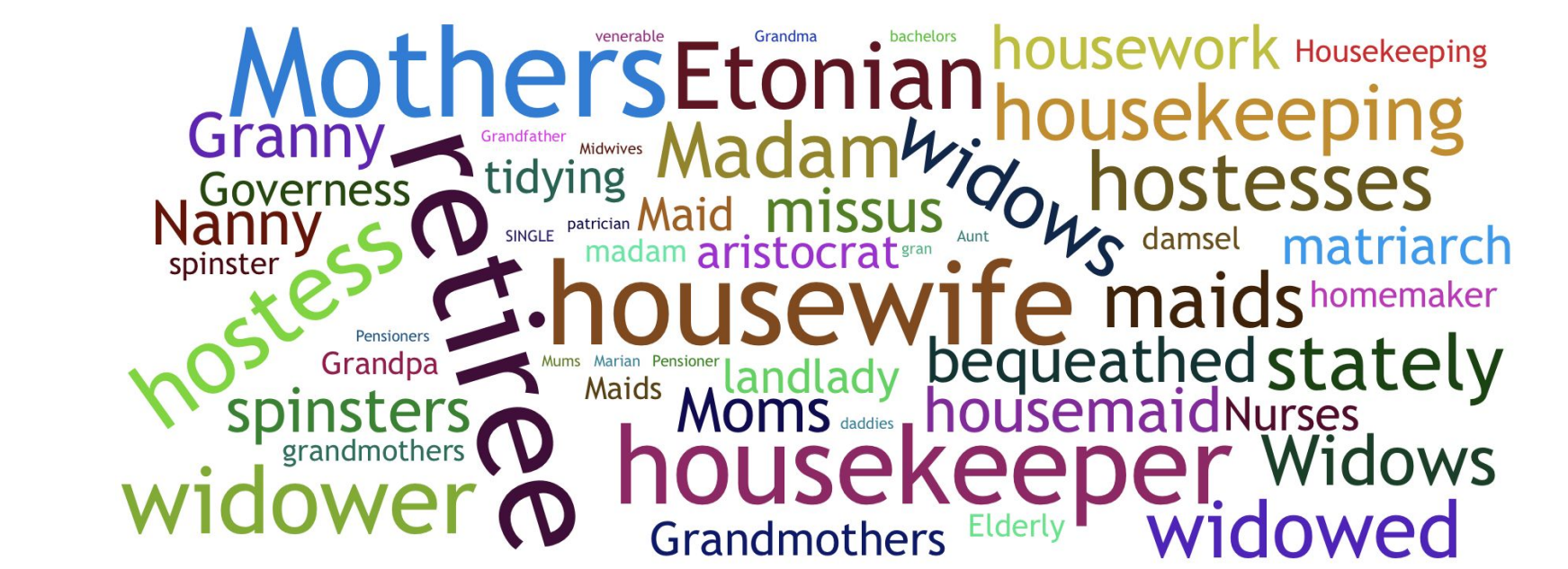The annotation task preserves the concept hierarchy.



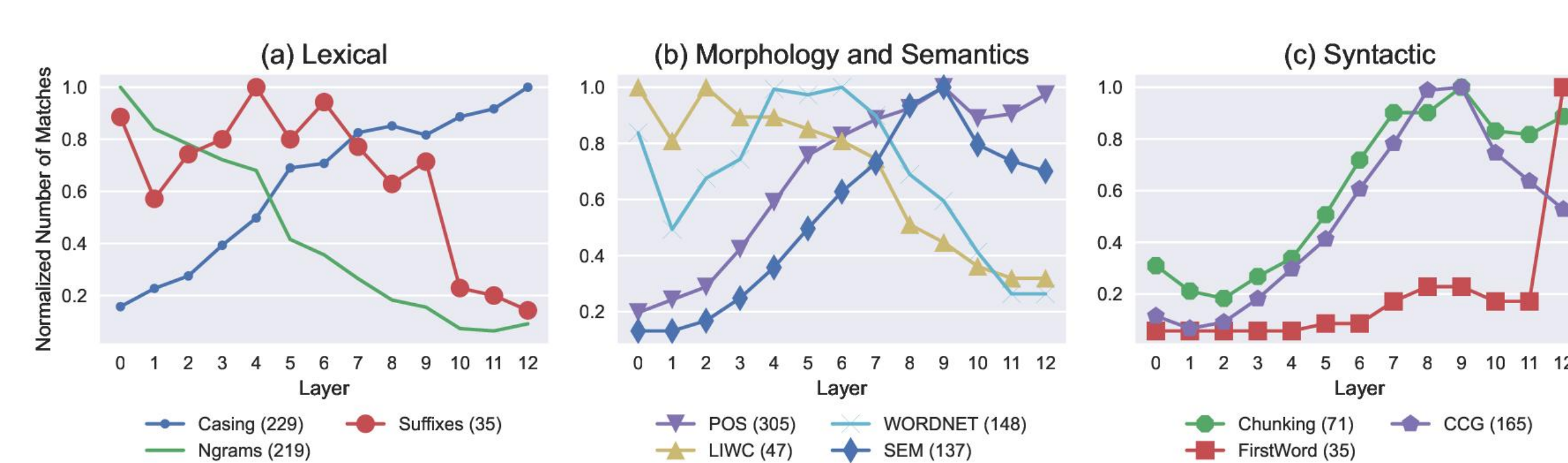The multifaceted nature captures diverse information.



## 6. Analysis

- Lexically similar but semantically different clusters based on the context
  - Decimal numbers that capture monetary values, e.g €9.6, $2.4M
  - Decimal numbers which appear as percentages, e.g. 9.6%, 2.4%

- Cluster shows potential biases present in the training data
  - Female roles such as mother, aunt, granny are grouped together with specific job roles such as housekeeper, maid and nanny



## 7. Alignment with pre-defined Concepts

- How much do latent BERT concepts align with pre-defined concepts ?
- Training data is annotated with pre-defined concepts
- A latent cluster is said to be aligned with the pre-defined concept if >=90% of its tokens belong to the pre-defined concept
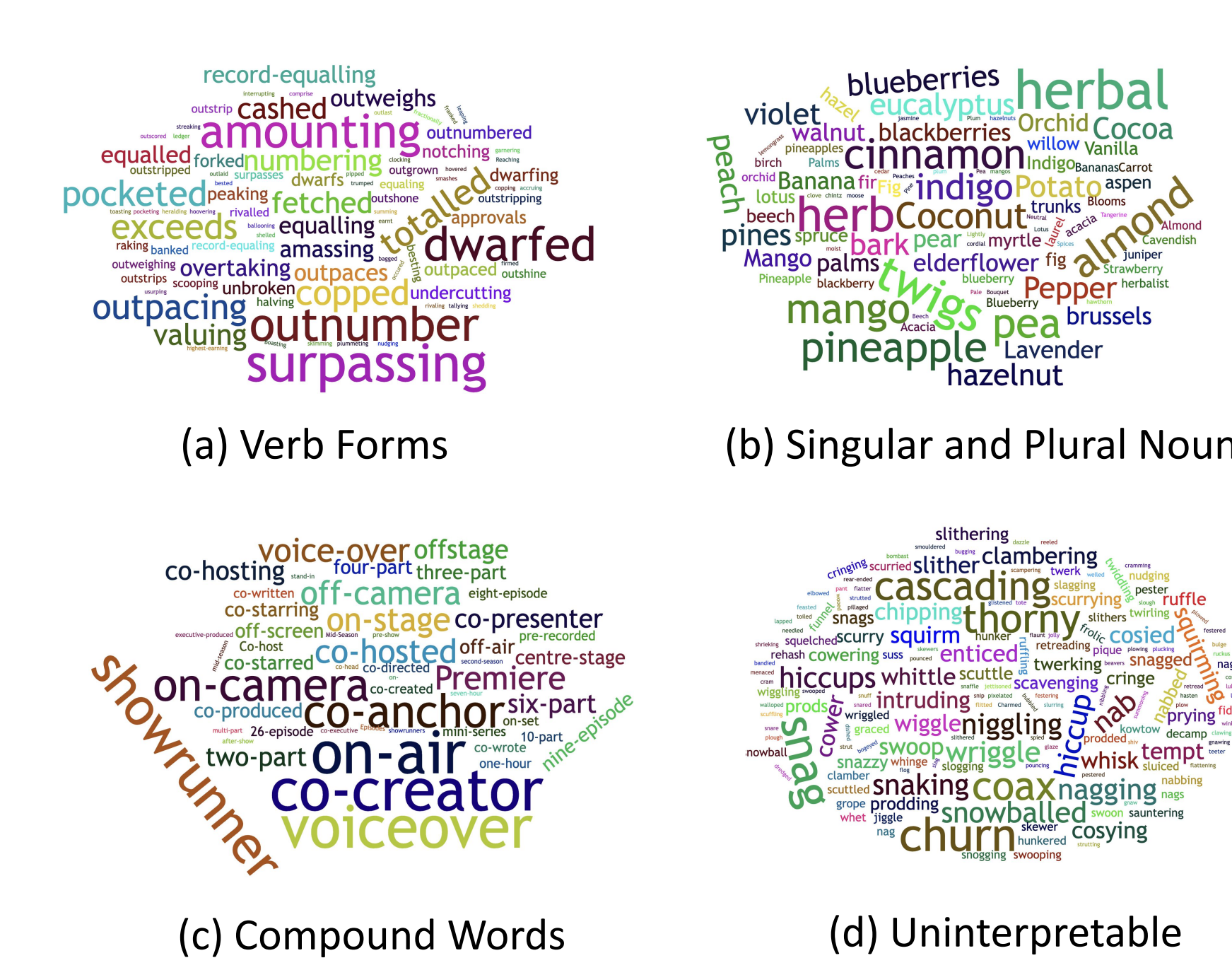
| Concepts Matches | Lexical | | | Morphology and Semantics | | | | Syntactic | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Ngram | Suffix | Casing | POS | SEM | LIWC | WordNet | CCG | Chunk | FW |
| | 20 (2.0%) | 5 (0.5%) | 229 (23%) | 297 (30%) | 96 (10%) | 15 (1.5%) | 39 (3.9%) | 87 (8.7%) | 63 (6.3%) | 35 (3.5%) |

Alignment of BERT concepts for layer 12 with pre-defined concepts

### Evolution of Concepts Across Layers



- Lower layers encode the lexical and meaning-related knowledge
- The encoded concepts evolve into representing linguistic hierarchy, in the higher layers, taking contextual information into account
- Higher alignment with lexical concepts (e.g. suffixes) in the lower layers
- Higher alignment with psycholinguistic concepts (e.g. LIWC) in the initial and middle layers
- Classical NLP concepts (e.g. POS, SEM, Chunking are captured in the middle and final layers

## 8. Unaligned Clusters

- What do the unaligned clusters represent?
- Compositionitonal clusters:
  - Figure (a) Verb forms  and (b) Singular/Plural Nouns
- Unaligned but explainable: Figure (c) Compound Words
- Uninterpretable Clusters: Figure (d) No meaningful relation



(a) Verb Forms

(b) Singular and Plural Nouns

(c) Compound Words

(d) Uninterpretable

## 9. BCN Dataset

To expand the manually annotated data:
- We trained a logistic classifier on the annotated concepts
- Predict the cluster id of new tokens from a large News data
- We only select a prediction when the classifier is 97% confident about its prediction
- **BCN consists of 174 concept labels and a total of 1M annotated instances**



**Application of BCN for Neuron Interpretation**
- Discover neurons learning a pre-defined concept
- BCN provides fine-grained concepts e.g. person names are split into finer categories based on geography
- Select a fine-grained concept; muslim names and a coarse concept person names
- Identify the neurons responsible for each concept

- We used Linguistic Correlation Analysis from the NeuroX toolkit to identify minimum number of neurons required for the concept
- We found only 19 neurons were required for the concept muslim names compared to 74 required for the concept person
- Therefore, this shows that BCN helps enable selection of specialized neurons responsible for very specific aspects of language