#### Holographic Reduced Representations

CS 135 - Brandeis University Fall 2024 James Pustejovsky

November 22, 2024

イロト イポト イヨト イヨト

#### Static and Contextualized Vector Composition

Holographic Reduced Representations (HRR)

Modeling GL Semantics with HRRs

#### Introduction

We explore three methods to represent sentences as vectors:

- Conventional vector composition
- Transformer-based contextualized embeddings
- Holographic Reduced Representations (HRRs)

We explore some semantic problems:

- Nouns as Vectors / Adjectives as Matrices
- Generative Lexicon encoded as Vector Binding
- Type Coercion in Vector Semantics

< ロ > < 同 > < 三 > < 三 >

#### What Are One-Hot Vectors?

- A one-hot vector is a binary vector used to represent categorical data.
- For a vocabulary of size V, each word is assigned a unique index i where:

$$\mathbf{v}_i = [0, 0, \dots, 1, \dots, 0]$$

イロト イポト イヨト イヨト

#### Example:

Vocabulary: ["cat", "dog", "fish"]
 "cat" → [1,0,0], "dog" → [0,1,0], "fish" → [0,0,1]

Limitations of One-Hot Vectors

High Dimensionality:

Dimension of vector = V

For large vocabularies (V > 100,000), the vectors become inefficient.

- Lack of Semantic Information:
  - No similarity between "cat" and "dog".
  - All vectors are orthogonal.
- Solution: Use Word2Vec to map one-hot vectors into dense, low-dimensional embeddings.

## From One-Hot Vectors to Word Embeddings (Word2Vec)

Word2Vec learns dense vector representations for words by analyzing their context in a corpus.

- Input: One-hot vector for each word.
- Output: Dense, low-dimensional embedding ( $\mathbf{w} \in \mathbb{R}^d$ ).

#### Key Idea: Distributional Hypothesis

Words that appear in similar contexts have similar meanings. Two training methods:

- Skip-gram: Predict context words from a target word.
- CBOW (Continuous Bag of Words): Predict a target word from context words.

イロト イポト イヨト イヨト

## Example: Word2Vec Conversion

- Vocabulary: ["cat", "dog", "fish"]
- One-hot vectors:

 $"\mathsf{cat}" = [1,0,0], \quad "\mathsf{dog}" = [0,1,0], \quad "\mathsf{fish}" = [0,0,1]$ 

Dense embeddings (Word2Vec output):

" cat" = [0.5, 0.1, 0.3], "dog" = [0.4, 0.2, 0.5], "fish" = [0.3, 0.8, 0.2]

イロト 不得 トイラト イラト 一日

These embeddings capture semantic similarity:

Similarity("cat", "dog") > Similarity("cat", "fish")

# What is Skip-gram?

- Word2Vec learns dense vector representations for words by predicting their context in a corpus.
- Skip-gram Model:
  - Predicts context words  $(w_c)$  given a target word  $(w_t)$ .
  - Objective: Maximize the probability of context words given the target word:

 $P(w_c|w_t)$ 

- Embedding Space:
  - ► Each word is mapped to a dense vector (d ~ 100 300).
  - Vectors capture semantic similarity (e.g., "king" and "queen").

イロト イヨト イヨト イヨト

## Details of the Skip-gram Model

#### **Objective Function**

For a given corpus, the Skip-gram model maximizes the conditional probability:

$$\prod_{t=1}^{T} \prod_{-c \leq j \leq c, j \neq 0} P(w_{t+j}|w_t)$$

イロト イヨト イヨト イヨト

where:

▶ *w<sub>t</sub>*: Target word.

•  $w_{t+i}$ : Context words within a window of size c.

# Details of the Skip-gram Model

#### Log-Likelihood

Taking the logarithm, the objective becomes:

$$\mathcal{L} = \sum_{t=1}^{T} \sum_{-c \leq j \leq c, j \neq 0} \log P(w_{t+j}|w_t)$$

#### Conditional Probability

The probability  $P(w_{t+j}|w_t)$  is modeled using softmax:

$$P(w_{t+j}|w_t) = \frac{\exp(\mathbf{u}_{w_{t+j}}^{\top}\mathbf{v}_{w_t})}{\sum_{w \in V} \exp(\mathbf{u}_w^{\top}\mathbf{v}_{w_t})}$$

イロト イヨト イヨト イヨト

**v**<sub>wt</sub>: Embedding for the target word.

**u**<sub>wt+i</sub>: Embedding for the context word.

# Skip-gram Training Steps

- 1. Initialize two embedding matrices:
  - W (target embeddings):  $|V| \times d$
  - W' (context embeddings):  $|V| \times d$
- 2. For each target word  $w_t$ :
  - Predict each context word  $w_c$  in the window [-c, c].
- 3. Compute the loss (negative log-likelihood):

$$\mathcal{L} = -\log P(w_c|w_t)$$

< ロ > < 同 > < 三 > < 三 >

4. Update W and W' using stochastic gradient descent (SGD).

## Computing Gradients for Skip-gram

For a single pair  $(w_t, w_c)$ , the loss is:

 $\mathcal{L} = -\log P(w_c|w_t)$ 

Gradient for Target Embedding  $(\mathbf{v}_{w_t})$ 

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}_{w_t}} = \mathbf{u}_{w_c} - \sum_{w \in V} P(w|w_t) \mathbf{u}_w$$

Gradient for Context Embedding  $(\mathbf{u}_{w_c})$ 

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}_{w_c}} = \mathbf{v}_{w_t} - \sum_{w \in V} P(w|w_t) \mathbf{v}_{w_t}$$

The gradients are used to update the embeddings via SGD.

# Negative Sampling: Reducing Computation

- Problem: Softmax requires summing over all words in the vocabulary (|V|).
- Solution: Use Negative Sampling to approximate softmax.

$$\log P(w_c|w_t) \approx \log \sigma(\mathbf{u}_{w_c}^{\top} \mathbf{v}_{w_t}) + \sum_{i=1}^{k} \mathbb{E}_{w \sim P_n(w)} \left[ \log \sigma(-\mathbf{u}_{w}^{\top} \mathbf{v}_{w_t}) \right]$$

- $P_n(w)$ : Noise distribution for negative samples.
- ▶ k: Number of negative samples per positive pair.

#### Advantages

- Reduces computation from O(|V|) to O(k).
- Focuses on distinguishing the target-context pairs from random noise.

# Worked Example: Skip-gram Training Step

Consider a toy vocabulary:  $V = {cat, dog, fish}$ .

- Target word:  $w_t = \text{cat.}$
- Context words:  $w_c \in \{ \text{dog}, \text{fish} \}$ .
- Embedding dimension: d = 2.
- Initialize embeddings:

$$\label{eq:v_cat} \textbf{v}_{cat} = [0.1, 0.3], \quad \textbf{u}_{dog} = [0.2, 0.4], \quad \textbf{u}_{fish} = [0.3, 0.1]$$

*,* T

Compute  $P(w_c|w_t)$  for  $w_c = \text{dog}$ :

$$P(w_c|w_t) = \frac{\exp(\mathbf{u}_{dog}^{\mathsf{v}}\mathbf{v}_{cat})}{\sum_{w \in V} \exp(\mathbf{u}_w^{\mathsf{T}}\mathbf{v}_{cat})}$$

Numerator:

$$\mathbf{J}_{dog}^{\top} \mathbf{v}_{cat} = (0.2)(0.1) + (0.4)(0.3) = 0.14$$

Denominator:

$$\mathsf{Sum} = \exp(0.14) + \exp(0.11) + \exp(0.03)$$

Result:

$$P(\text{dog}|\text{cat}) = \frac{\exp(0.14)}{\exp(0.14) + \exp(0.11) + \exp(0.03)} = \dots$$

## Final Word Embeddings

- After training, each word has two embeddings:
  - $\blacktriangleright$  **v**<sub>w<sub>t</sub></sub>: Represents the word as a target.
  - $\mathbf{u}_{w_c}$ : Represents the word as a context.
- Combine these embeddings (e.g., by averaging) to create the final word vector:

$$\mathbf{w} = \frac{\mathbf{v}_{w_t} + \mathbf{u}_{w_c}}{2}$$

A (10) × (10) × (10) ×

These embeddings capture semantic similarity and are used in downstream tasks.

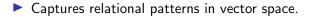
# Using Word2Vec for Analogical Reasoning

 Analogical reasoning involves finding relationships between pairs of words or concepts.

"king is to queen as man is to woman."

Represented mathematically as:

king – man  $\approx$  queen – woman.



# Latent Semantic Analysis (LSA)

- LSA uses singular value decomposition (SVD) to reduce the dimensionality of term-document matrices.
- Represents words and documents as vectors in a semantic space:

$$\mathbf{M} \in \mathbb{R}^{|V| \times |D|} \quad \rightarrow \quad \mathbf{M}_{\mathsf{reduced}} \in \mathbb{R}^{|V| \times k}$$

Captures semantic relationships:

cosine similarity between word vectors reflects semantic similarity.

< ロ > < 同 > < 回 > < 回 >

#### Limitations

- LSA focuses on co-occurrence, not relational patterns.
- Cannot explicitly represent analogies.

## The Parallelogram Hypothesis

- Hypothesis: Analogical reasoning can be represented as geometric relationships in vector space.
- Example:

```
king – man + woman \approx queen.
```

- Geometric Interpretation:
  - The vector from man to king is parallel to the vector from woman to queen.

 Visualized as a parallelogram: Given three points, solve for the fourth:

#### queen = king - man + woman.

#### How Word2Vec Derives Analogies

- Word2Vec learns dense word embeddings that capture semantic and syntactic relationships.
- Relationships are encoded in the directions between vectors.
- Analogy-solving formula:

$$\mathbf{w}_4 = \arg \max_{\mathbf{w} \in V} \cos \left( \mathbf{w}, \mathbf{w}_2 - \mathbf{w}_1 + \mathbf{w}_3 \right),$$

・ロ・ ・ 回 ・ ・ ヨ ・ ・ ヨ ・ ・

where:

w<sub>1</sub> = man, w<sub>2</sub> = king, w<sub>3</sub> = woman,
 w<sub>4</sub> = queen.

## Worked Example: Word2Vec Analogy

Example: Solve "king is to queen as man is to woman":

Vectors:

king = [0.8, 0.6], queen = [0.9, 0.7],

$$man = [0.2, 0.4], woman = [0.3, 0.5].$$

Compute:

Result:

queen = king 
$$-$$
 man  $+$  woman.  
queen =  $[0.8, 0.6] - [0.2, 0.4] + [0.3, 0.5].$ 

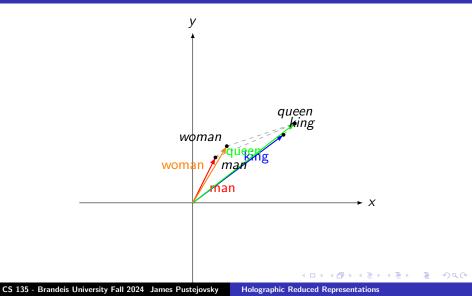
queen = 
$$[0.9, 0.7]$$
.

イロト イヨト イヨト イヨト

Static and Contextualized Vector Composition

Holographic Reduced Representations (HRR) Modeling GL Semantics with HRRs

#### Vector Analogy with Parallelogram Visualization



## Why Does Word2Vec Work?

- Co-occurrence Modeling:
  - Word2Vec captures context relationships via training on skip-grams.
- Semantic Directionality:
  - Embeddings encode directional relationships (e.g., gender, tense).
- Vector Arithmetic:
  - The geometry of word embeddings allows analogical reasoning through addition and subtraction.

< ロ > < 同 > < 三 > < 三 >

## What Word2Vec Does Not Explain

- Syntax-Semantics Interface:
  - Analogies focus on semantics; no explicit representation of syntactic structure.
- Complex Analogies:
  - Cannot handle multi-step or hierarchical relationships.
- Context Dependence:
  - Word2Vec embeddings are static, ignoring polysemy and contextual nuances.
- Empirical Limitations:
  - Only works well for analogies seen in training or closely related domains.

< ロ > < 同 > < 三 > < 三 >

# Key Takeaways

- Analogical reasoning is a fundamental capability of word embeddings like Word2Vec.
- The parallelogram hypothesis explains how analogies are geometrically encoded in vector space.
- Limitations:
  - Word2Vec does not capture syntax or hierarchical relationships.

- Contextualized embeddings (e.g., BERT) address some limitations but are less interpretable.
- Analogical reasoning with vectors demonstrates the power and constraints of distributional semantics.

#### Transformer-Based Contextualized Embeddings

Overview of Self-Attention in Transformers

- Self-attention computes relationships between tokens in a sentence.
- Outputs contextualized representations for each token.

#### Self-Attention Formula

For query (Q), key (K), and value (V) matrices:

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- ▶  $Q, K, V \in \mathbb{R}^{n \times d_k}$ , where *n* is the number of tokens and *d\_k* is the embedding size.
- Each token generates its own query, key, and value vectors.

#### Step 1: Input to Embeddings

Given a sentence:

"The bank will not approve the loan."

► Tokens: [The, bank, will, not, approve, the, loan].

• Embedding dimension:  $d_k = 4$  (for simplicity).

Token embeddings (random initialization for this example):

$$\textbf{x}_1 = [1,0,1,0], \, \textbf{x}_2 = [0,1,0,1], \dots$$

イロト イヨト イヨト イヨト

#### Step 2: Compute Query, Key, and Value Matrices

Each token embedding is projected into query, key, and value spaces using learned weight matrices:

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V$$

Example weights  $(W_Q, W_K, W_V \in \mathbb{R}^{d_k \times d_k})$ :

$$W_Q = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix}, \ W_K = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \end{bmatrix}, \ \dots$$

For token 1  $(x_1 = [1, 0, 1, 0])$ :

$$Q_1 = \mathbf{x}_1 W_Q = \begin{bmatrix} 1, 0, 1, 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 2, 0, 1, 0 \end{bmatrix}$$

## Step 3: Compute Attention Scores

Compute scaled dot-product attention:

$$\mathsf{Attention}(Q,K,V) = \mathsf{softmax}\left(\frac{QK^{\mathsf{T}}}{\sqrt{d_k}}\right)V$$

For  $Q_1$  and  $K_2$ :

$$Q_1 \cdot K_2 = [2, 0, 1, 0] \cdot [0, 1, 0, 1] = 0$$

Attention scores matrix:

$$Scores_{i,j} = \frac{Q_i \cdot K_j}{\sqrt{d_k}}, \quad i,j \in \{1,\ldots,n\}$$

Normalize scores using softmax:

$$\mathsf{softmax}(z_i) = rac{\mathsf{exp}(z_i)}{\sum_j \mathsf{exp}(z_j)}$$

・ 同 ト ・ ヨ ト ・ ヨ ト

### Step 4: Compute Weighted Values

For token 1:

Attention
$$(Q_1, K, V) = \operatorname{softmax}\left(\frac{Q_1K^T}{\sqrt{d_k}}\right)V$$

Example:

$$Scores_{1,} = softmax\left(\frac{[0,1,2]}{\sqrt{4}}\right) = [0.04, 0.11, 0.85]$$

Use scores to compute weighted sum:

$$\mathbf{z}_1 = \sum_{j=1}^n \operatorname{Scores}_{1,j} \cdot V_j$$

#### Worked Example: Final Contextualized Embeddings

For each token *i*, compute:

$$\mathbf{z}_i = \sum_{j=1}^n \operatorname{Attention}(Q_i, K_j, V_j)$$

Result:

$$Z = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \mathbf{z}_n \end{bmatrix}, \quad Z \in \mathbb{R}^{n \times d_k}$$

Aggregated sentence embedding:

$$\mathbf{S}_{\mathsf{Transformer}} = \mathsf{Mean}(Z)$$

or use special token [CLS].

CS 135 - Brandeis University Fall 2024 James Pustejovsky

(4月) トイヨト イヨト

# Motivation for Hypervectors and Hyperdimensional Computing

- ▶ Hypervectors: High-dimensional vectors ( $d \gg 1000$ ) used to represent information in a distributed manner.
- Inspired by the properties of the brain:
  - Robustness to noise.
  - Ability to store and retrieve large amounts of information.

#### Key idea:

Complex structures can be represented as combinations of simple high-dimensional vectors.

イロト イポト イヨト イヨト

#### What is a Hyperdimensional Vector Space?

- A hyperdimensional vector space is a high-dimensional space (d ≫ 1000) used to represent information.
- Hypervectors ( $\mathbf{v} \in \mathbb{R}^d$ ):
  - Randomly initialized.
  - High dimensionality ensures approximate orthogonality between vectors.

#### Properties of High-Dimensional Spaces

Orthogonality:

For random vectors  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ :  $\mathbf{v} \cdot \mathbf{w} \approx 0$  (if  $\mathbf{v} \neq \mathbf{w}$ ).

Stability:

$$\|\mathbf{v}_1 + \mathbf{v}_2 + \cdots + \mathbf{v}_n\|$$
 grows with  $n$ .

#### Capacity:

Large spaces can encode exponentially many distinct patterns.

(本部) ( 문) ( 문) ( 문

## How do Hyperdimensional Spaces Encode Patterns"

Key Property: In a hyperdimensional space ℝ<sup>d</sup> with d ≫ 1000:

Number of distinct patterns is exponential in d.

- Intuition:
  - A random hypervector  $\mathbf{v} \in \mathbb{R}^d$  has d components.
  - ► Each component can take on many possible values (e.g., [-1,1] for bipolar vectors or ℝ for real-valued vectors).
- Mathematical Argument:
  - ▶ Consider *d*-dimensional binary vectors {0,1}<sup>*d*</sup>:

Total number of distinct vectors:  $2^d$ .

For real-valued or bipolar vectors, the number of distinct patterns grows even faster.

## Geometric Perspective: Orthogonality in High Dimensions

High-dimensional spaces have the property that random vectors are nearly orthogonal:

 $\textbf{v}_1 \cdot \textbf{v}_2 \approx 0, \quad \text{if } \textbf{v}_1, \textbf{v}_2 \text{ are random}.$ 

#### Implication:

- You can generate exponentially many random hypervectors that are distinguishable (linearly independent or approximately orthogonal).
- Example:
  - ▶ In  $\mathbb{R}^{10,000}$ , billions of random vectors will have dot products close to zero.

# Superposition and Binding in Hyperdimensional Spaces

Superposition (addition):

$$\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2 + \cdots + \mathbf{v}_n.$$

- ► Even with n ≫ 1, the resulting vector is still distinguishable due to high dimensionality.
- Binding (e.g., circular convolution):

$$\mathbf{c} = \mathbf{a} \circledast \mathbf{b}.$$

- Each binding operation produces a new, distinguishable pattern.
- Exponentially Growing Combinations:
  - *n* hypervectors can generate:
    - $2^n$  combinations via binding and superposition.

イロト イヨト イヨト

## Worked Example: Exponential Growth of Patterns

Consider d = 10,000 and binary vectors  $\{0,1\}^d$ :

Total possible distinct vectors:

2<sup>10,000</sup> (an astronomically large number).

Now allow for superposition and binding:

- Superposition combines *n* vectors into a unique vector.
- Binding generates entirely new patterns:

 $\mathbf{a} \circledast \mathbf{b}$  is unique for any  $\mathbf{a}, \mathbf{b}$ .

イロト イヨト イヨト イヨト

Result:

With high-dimensional vectors, you can encode exponentially many relationships.

# Key Benefits of Exponentially Large Spaces

#### Robustness:

- Small errors (noise) in the components of hypervectors do not significantly affect overall distinguishability.
- Scalability:
  - Exponentially large capacity ensures scalability for encoding large vocabularies, complex patterns, and relationships.

#### Expressiveness:

- Binding and superposition operations allow for compositional representations (e.g., hierarchical structures or analogies).
- Similarity Preservation:
  - High-dimensional vectors can preserve similarity in the space (e.g., similar words have closer embeddings).

#### Comparison: One-Hot Vectors vs. Hypervectors

Dimensionality:

- One-hot: |V| (grows with vocabulary size).
- Hypervectors: d (fixed, large dimensionality, e.g., d = 10,000).
- Orthogonality:

One-hot:

$$\mathbf{v}_i \cdot \mathbf{v}_j = 0$$
 for  $i \neq j$ .

Hypervectors:

 $\mathbf{v}_i \cdot \mathbf{v}_j \approx 0$  (approximate for random vectors).

イロト イポト イヨト イヨト

#### Representation Power:

One-hot: Only encodes identity.

Hypervectors: Encodes identity, similarity, and relationships.

### **Computational Distinctions**

- Storage Requirements:
  - One-hot: Requires a vector of size |V| for each token.
  - Hypervectors: Fixed size d, independent of vocabulary size.
- Operations:
  - One-hot: No meaningful operations (e.g., addition, multiplication).
  - Hypervectors: Supports binding, superposition, and correlation.

$$\mathbf{a} \circledast \mathbf{b}, \quad \mathbf{a} + \mathbf{b}, \quad \mathbf{a} \circledast \mathbf{b}^{-1}.$$

< ロ > < 同 > < 三 > < 三 >

Scalability:

- One-hot: Becomes infeasible for large vocabularies (|V| ≫ 10<sup>6</sup>).
- Hypervectors: Efficient for large vocabularies due to fixed dimensionality.

# Advantages of Hypervector Encodings

Compositionality:

Represent relationships through binding:

relation =  $a \circledast b$ .

Combine multiple pieces of information:

 $context = v_1 + v_2 + v_3.$ 

Noise Tolerance:

Small changes to components do not disrupt overall encoding.

- Similarity Preservation:
  - Similar inputs produce similar hypervectors, enabling clustering and matching.
- Scalability:
  - Fixed-dimensional encoding handles large vocabularies and complex structures.

#### Comparison: One-Hot Vectors vs. Hypervectors

Feature	One-Hot Vectors	Hypervectors
Dimensionality	V  (vocab size)	Fixed <i>d</i> (e.g., 10,000)
Orthogonality	Exact	Approximate
Representation	Identity only	Identity $+$ relationships
Operations	None	Binding, superposition, correlation
Storage	Large for large $ V $	Fixed-size
Scalability	Limited	High

イロト イヨト イヨト イヨト

## Why Orthonormality is Important

In high-dimensional spaces, random hypervectors are approximately orthonormal:

$$\mathbf{v} \cdot \mathbf{w} pprox \mathbf{0}, \quad \|\mathbf{v}\| = \|\mathbf{w}\| = 1.$$

 Key Implication: Vectors do not interfere with each other in superposition or binding.

 Superposition: Combines multiple vectors while keeping them distinguishable.

$$\mathbf{s} = \mathbf{v}_1 + \mathbf{v}_2 + \cdots + \mathbf{v}_n$$

 Binding: Combines vectors into unique encodings using circular convolution.

$$\mathbf{b}=\mathbf{v}_1 \circledast \mathbf{v}_2$$

< ロ > < 同 > < 三 > < 三 >

## Circular Convolution: Binding Vectors

Circular convolution is defined as:

$$(\mathbf{a} \circledast \mathbf{b})_i = \sum_{j=0}^{d-1} a_j \cdot b_{(i-j) \mod d}.$$

#### Properties of Circular Convolution

$$\mathsf{a} \circledast \mathsf{b} \in \mathbb{R}^d$$

- Uniqueness: Produces a distinct vector for each pair of inputs.
- Approximate Inverse:

$$\mathbf{a} \circledast \mathbf{b} \circledast \mathbf{b}^{-1} \approx \mathbf{a}$$

A (10) × (10) × (10) ×

## Correlation: Unbinding Vectors

Circular correlation retrieves one vector from a bound pair:

$$(\mathbf{a} \circledast \mathbf{b}) \circledast \mathbf{b}^{-1} \approx \mathbf{a}.$$

#### Definition of Circular Correlation

Circular correlation is defined as:

$$(\mathbf{c} \circledast \mathbf{b}^{-1})_i = \sum_{j=0}^{d-1} c_j \cdot b_{(j-i) \mod d}.$$

#### Key Insights

- Uses the approximate orthonormality of random vectors.
- Recovers the original vector when the bound pair is unbound.

### Approximate Orthogonality in High Dimensions

For two random hypervectors  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ :

$$\mathbf{v}\cdot\mathbf{w}=\sum_{i=1}^d v_iw_i.$$

If  $\mathbf{v}_i$  and  $\mathbf{w}_i$  are independent and zero-mean:

Expected value:

$$\mathbb{E}[\mathbf{v}\cdot\mathbf{w}]=0.$$

Variance decreases with dimensionality:

$$\mathsf{Var}[\mathbf{v}\cdot\mathbf{w}]=O\left(rac{1}{d}
ight).$$

イロト イポト イヨト イヨト

For large d, the dot product is negligibly small:

 $\mathbf{v} \cdot \mathbf{w} \approx 0.$ 

## Implications for Encoding in NLP

 Superposition: Adding hypervectors preserves distinguishability:

$$\mathbf{s} = \mathbf{v}_1 + \mathbf{v}_2 \quad \Rightarrow \quad \mathbf{s} \cdot \mathbf{v}_1 \gg \mathbf{0}.$$

Binding: Convolution produces unique encodings:

$$\mathbf{b}=\mathbf{v}_1 \circledast \mathbf{v}_2.$$

Since  $\textbf{v}_1\cdot \textbf{v}_2\approx 0,$  the result is not confounded by interference.

Unbinding: Correlation retrieves components reliably:

$$(\mathbf{b} \circledast \mathbf{v}_2^{-1}) \approx \mathbf{v}_1.$$

< ロ > < 同 > < 三 > < 三 >

# Example: Semantic Role Binding

Represent the sentence "The dog chased the ball":

- ► Words: **dog**, **chased**, **ball**.
- Roles: subject, verb, object.

Encoding:

$$S = (dog \circledast subject) + (chased \circledast verb) + (ball \circledast object).$$

Retrieval:

Retrieve the subject:

$$\operatorname{dog} \approx \mathbf{S} \circledast \operatorname{subject}^{-1}$$
.

A (10) × (10) × (10) ×

# Example: Sequential Encoding

Encode "The dog sleeps":

- ► Words: the, dog, sleeps.
- Positional encoding:

sequence = the + 
$$Perm(dog) + Perm^{2}(sleeps)$$
.

Retrieval:

Retrieve "dog" by reversing the permutation:

 $\log \approx$  sequence  $\circledast$  Perm<sup>-1</sup>.

イロト イヨト イヨト イヨト

# Vector Algebraic Composition Operations

Three primary operations are used in hyperdimensional computing:

Superposition (Addition):

$$\mathbf{c} = \mathbf{a} + \mathbf{b}$$

Combines vectors while preserving their individual contributions.

Binding (Multiplication or Convolution):

 $\mathbf{c} = \mathbf{a} \circledast \mathbf{b}$ 

Creates a unique composite vector that is distinct from the inputs.

Permutation:

Random reordering of vector components to represent positional information.

# Convolution as Binding

#### Circular Convolution

The binding operation in HRRs is defined as circular convolution:

$$c_i = \sum_{j=0}^{d-1} a_j \cdot b_{(i-j) \mod d}$$

Here:

• 
$$\mathbf{a} = [a_0, a_1, \dots, a_{d-1}]$$

**b** = 
$$[b_0, b_1, \dots, b_{d-1}]$$

- Properties:
  - Produces a vector of the same dimension d.
  - Distributes information of a and b across c.
  - Approximately invertible.

### Worked Example: Circular Convolution

#### Given:

$$\mathbf{a} = [1, 2, 3], \quad \mathbf{b} = [0, 1, 0]$$

Compute:

$$c_0 = a_0 \cdot b_0 + a_1 \cdot b_2 + a_2 \cdot b_1 = 1 \cdot 0 + 2 \cdot 0 + 3 \cdot 1 = 3$$

$$c_1 = a_0 \cdot b_1 + a_1 \cdot b_0 + a_2 \cdot b_2 = 1 \cdot 1 + 2 \cdot 0 + 3 \cdot 0 = 1$$

$$c_2 = a_0 \cdot b_2 + a_1 \cdot b_1 + a_2 \cdot b_0 = 1 \cdot 0 + 2 \cdot 1 + 3 \cdot 0 = 2$$
Result:

$$c = [3, 1, 2]$$

イロト イヨト イヨト イヨト

臣

#### Deconvolution as Unbinding

#### Deconvolution

To retrieve **a** from **c** and **b**, perform circular correlation:

$$a_i = \sum_{j=0}^{d-1} c_j \cdot b_{(i-j) \mod d}$$

Inverse Property:

$$\mathbf{a}\approx\mathbf{c}\circledast\mathbf{b}^{-1}$$

< ロ > < 同 > < 三 > < 三 >

Enables retrieval of the original components bound together.

## Summary of HRR Operations

Binding (Encoding):

$$\mathbf{binding} = \mathbf{w}_{\mathsf{word}} \circledast \mathbf{r}_{\mathsf{role}}$$

Superposition:

$$\mathbf{S}_{\mathsf{HRR}} = \mathbf{S} + \mathbf{V} + \mathbf{Neg} + \mathbf{O}$$

Unbinding (Decoding):

$$\mathsf{w}_{\mathsf{word}} pprox \mathbf{S}_{\mathsf{HRR}} \circledast \mathsf{r}_{\mathsf{role}}^{-1}$$

(1) マント (1) マント (1) マント

HRRs provide a robust framework for representing and manipulating structured information in high-dimensional spaces.

#### Comparison of Methods

#### A comparison of the three methods:

Feature	Transformer	Conventional	HRR
Context-Sensitivity	High	None	Moderate
Syntactic Structure	Implicit	Ignored	Explicit
Polysemy Handling	Excellent	Poor	Limited
Invertibility	No	No	Yes

イロト イヨト イヨト イヨト

臣

## Overview of Baroni & Zamparelli's Theory

Nouns: Represented as dense vectors.

$$\mathbf{n} \in \mathbb{R}^d$$

Adjectives: Represented as linear transformations (matrices).

$$\mathbf{A} \in \mathbb{R}^{d imes d}$$

 Modification: Apply the adjective to the noun using matrix-vector multiplication.

Modified noun:  $\mathbf{n}' = \mathbf{A} \cdot \mathbf{n}$ 

< ロ > < 同 > < 三 > < 三 >

## Worked Example: Baroni & Zamparelli's Theory

Let:

• Noun: 
$$\mathbf{n} = [1, 0, 1]^{T}$$
  
• Adjective:  $\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$ 

Compute:

$$\mathbf{n}' = \mathbf{A} \cdot \mathbf{n} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Interpretation:

- Adjective transforms the noun's semantic space.
- Matrix captures how adjectives modify meaning (e.g., "red" or "big").

## Strengths and Limitations of Baroni & Zamparelli's Theory

#### Strengths:

Captures compositional semantics via linear transformations.

Allows for a systematic representation of adjective effects.

#### Limitations:

- High parameter cost ( $d^2$  parameters per adjective).
- Limited interpretability of learned matrices.
- Ignores distributed binding (no explicit roles or structure).

(4回) (4回) (4回)

### Overview of HRR Approach

► Nouns: Represented as high-dimensional hypervectors.

$$\mathbf{n} \in \mathbb{R}^d, \quad d \gg 1000$$

Adjectives: Represented as hypervectors.

$$\mathbf{a} \in \mathbb{R}^d$$

Binding: Adjectives bind to nouns using circular convolution.

$$\mathbf{n}' = \mathbf{a} \circledast \mathbf{n}$$

Superposition: Combine multiple adjective-noun pairs.

$$\bm{S} = \bm{n}_1' + \bm{n}_2' + \cdots$$

# Worked Example: HRR Approach

#### Let:

- Noun hypervector:  $\mathbf{n} = [1, 0, 1]$
- Adjective hypervector:  $\mathbf{a} = [0, 1, 0]$

Compute binding via circular convolution:

$$c_0 = a_0 \cdot n_0 + a_1 \cdot n_2 + a_2 \cdot n_1 = 0 \cdot 1 + 1 \cdot 1 + 0 \cdot 0 = 1$$
  

$$c_1 = a_0 \cdot n_1 + a_1 \cdot n_0 + a_2 \cdot n_2 = 0 \cdot 0 + 1 \cdot 1 + 0 \cdot 1 = 1$$
  

$$c_2 = a_0 \cdot n_2 + a_1 \cdot n_1 + a_2 \cdot n_0 = 0 \cdot 1 + 1 \cdot 0 + 0 \cdot 1 = 0$$

Resulting vector:

$$\boldsymbol{n}' = [1,1,0]$$

(4月) トイヨト イヨト

### Comparison: Baroni & Zamparelli vs. HRR

Feature	Baroni & Zamparelli	HRR
Representation of Adjectives	Matrices $(d \times d)$	Hypervectors (d)
Composition Operation	Matrix-Vector Multiplication	Circular Convolution
Dimensionality	$O(d^2)$ (scaling issue)	Fixed (d)
Interpretability	Low	Moderate
Flexibility (e.g., roles)	Limited	High (binding and superposition)

イロト イヨト イヨト イヨト

臣

## The Role of Telic in Generative Lexicon

The Telic role in GL Theory captures the purpose or function of an entity.

Telic(pen) = write-with, Telic(car) = drive

- Adjectives modify nouns by binding to specific Qualia roles, including the Telic role.
- Disambiguation occurs when an adjective aligns with the Telic role of the noun.

#### Examples of Telic-Driven Disambiguation

- ▶ "Fast car"  $\rightarrow$  Telic: drive (interpreted as speed when driving).
- "Good pen"  $\rightarrow$  Telic: write-with (interpreted as quality in writing).
- ► "Loud speaker" → Telic: produce-sound (interpreted as volume of sound production).

# Selective Binding with Telic Role

Adjective-noun composition involves selective binding:

```
Adjective 

    Telic(Noun)
```

- This highlights the Telic role of the noun as the locus of modification.
- Formalization:

 $\mathbf{a} \circledast \mathbf{r}_{\mathsf{Telic}} \circledast \mathbf{n}$ 

#### Interpretation of Adjective Modification

The adjective binds to the Telic role, influencing how the noun is interpreted in context.

"Fast car" = fast 
$$\circledast$$
 drive(car)

イロト イポト イヨト イヨト

# HRR Representation of Telic Role

- Nouns: Represented as hypervectors (n).
- Adjectives: Represented as hypervectors (a).
- ► Telic Role: Represented as a hypervector (**r**<sub>Telic</sub>).
- Binding: Use circular convolution to encode adjective modification of the Telic role.

 $\textbf{composition} = a \circledast r_{\textsf{Telic}} \circledast n$ 

< ロ > < 同 > < 三 > < 三 >

 Superposition: Combine multiple adjective-noun pairs for broader contexts.

# Worked Example: "Fast Car"

Given:

• Noun: "car" 
$$(\mathbf{n} = [1, 0, 1])$$
.

• Telic Role: "drive" (
$$\mathbf{r}_{\text{Telic}} = [1, 1, 0]$$
).

Compute:

$$binding = a \circledast r_{Telic} \circledast n$$

Step 1 (Adjective-Telic binding):

$$c_0 = a_0 r_0 + a_1 r_2 + a_2 r_1 = 0 \cdot 1 + 1 \cdot 0 + 0 \cdot 1 = 0$$

Step 2 (Adjective-Telic-Noun binding):

$$c_0' = c_0 n_0 + c_1 n_2 + c_2 n_1 = \dots$$

Result:

#### $\textbf{binding} = \dots$

The final vector encodes "fast car" in terms of its Telic role, =>> = ∽ << CS 135 - Brandeis University Fall 2024 James Pustejovsky Holographic Reduced Representations

# Worked Example: "Good Pen"

#### Given:

- ► Adjective: "good" (**a** = [1,0,1]).
- Telic Role: "write-with" ( $\mathbf{r}_{\text{Telic}} = [0, 1, 1]$ ).

Compute:

$$binding = a \circledast r_{\mathsf{Telic}} \circledast n$$

Result:

#### binding $= \dots$

イロト イポト イヨト イヨト

This vector captures the modification of "pen" by "good" with respect to the Telic role, emphasizing its quality in writing.

# What is Type Coercion?

- Type Coercion occurs when a verb's argument requires a type mismatch to be resolved.
- The mismatched argument is coerced into the required type using its Qualia Structure.

#### Example: "Mary enjoyed a coffee"

- Verb: "enjoy" requires an event as its object.
- Noun: "a coffee" is an entity, not an event.
- Coercion: The Telic role of "coffee" provides the event "drink a coffee".

イロト イポト イヨト イヨト

"Mary enjoyed a coffee"  $\Rightarrow$  "Mary enjoyed drinking a coffee"

# Type Coercion Using Qualia Structure

- ► The Qualia Structure of "coffee":
  - Formal: beverage.
  - Constitutive: made-of-water.
  - Telic: drink.
  - Agentive: *brewed*.

► The Telic role provides the required event for coercion:

Telic(coffee) = drink(coffee)

#### Coercion Process

1. Verb identifies a type mismatch (*entity* vs. *event*). 2. Use the Qualia Structure to resolve the mismatch. 3. Bind the Telic role to the object and recompose:

< ロ > < 同 > < 三 > < 三 >

## HRR Representation of Coercion

- Nouns: Represented as hypervectors (n).
- ► Telic Role: Represented as a hypervector (**r**<sub>Telic</sub>).
- ▶ Verb: Represented as a hypervector (**v**<sub>enjoy</sub>).
- Coercion: Bind the verb to the Telic role of the object:

 $composition = v_{enjoy} \circledast r_{Telic} \circledast n_{coffee}$ 

# Worked Example: "Mary Enjoyed a Coffee"

Given:

▶ Noun hypervector: 
$$\mathbf{n}_{coffee} = [1, 0, 1].$$

• Telic Role: 
$$\mathbf{r}_{\text{Telic}} = [0, 1, 0].$$

• Verb: 
$$\mathbf{v}_{enjoy} = [1, 1, 0].$$

Compute coercion:

 $\textbf{composition} = v_{\text{enjoy}} \circledast r_{\text{Telic}} \circledast n_{\text{coffee}}$ 

Step 1 (Verb-Telic binding):

$$c_0 = v_0 r_0 + v_1 r_2 + v_2 r_1 = 1 \cdot 0 + 1 \cdot 0 + 0 \cdot 1 = 0$$

Step 2 (Verb-Telic-Noun binding):

$$c_0' = c_0 n_0 + c_1 n_2 + c_2 n_1 = \dots$$

Result:

$$\textbf{composition} = [\dots]$$

This vector encodes the coerced meaning "enjoy drinking a coffee."

### Denotational Semantics and Type-Theoretic Interpretation

- $\blacktriangleright \ [\![expression]\!] \rightarrow mathematical \ structure$
- Examples:
  - Words:  $\llbracket dog \rrbracket = vector in \mathbb{R}^d$ .
  - Sentences:  $\llbracket$  Dogs bark"  $\rrbracket$  = truth value {0, 1}.
- In vector-based semantics:
  - Nouns:  $[\![N]\!] : \mathbb{R}^d$
  - Adjectives:  $[Adj] : \mathbb{R}^d \to \mathbb{R}^d$
  - Verbs:  $\llbracket V \rrbracket : \mathbb{R}^d \to \mathbb{R}^d$
- Sentence composition:

 $[\![" Dogs bark"]\!] =$  function application or combination of vectors.

イロト イポト イヨト イヨト

## Conventional Vector Composition

- Words are vectors:  $\llbracket w \rrbracket \in \mathbb{R}^d$ .
- Composition is performed using vector addition or pointwise multiplication.

#### **Denotational Semantics**

Let 
$$w_1, w_2 \in \mathbb{R}^d$$
, then:

$$\llbracket " \, \mathsf{fast} \, \, \mathsf{car}" \, 
rbracket = \mathbf{v}_\mathsf{fast} + \mathbf{v}_\mathsf{car}.$$

Type-Theoretic Interpretation

▶ Nouns:  $e : \mathbb{R}^d$ 

• Adjectives: 
$$e \rightarrow e : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Composition:

$$\mathsf{Adj}(\mathsf{N}): \mathbb{R}^d o \mathbb{R}^d$$

$$\llbracket "\operatorname{fast"} \rrbracket (\llbracket "\operatorname{car"} \rrbracket) = \mathbf{v}_{\operatorname{fast}} + \mathbf{v}_{\operatorname{car}}.$$

### Baroni and Zamparelli: Adjectives as Matrices

Words are assigned different types:

▶ Nouns:  $\mathbb{R}^d$ 

• Adjectives:  $\mathbb{R}^d \to \mathbb{R}^d$  (matrices)

► Composition uses matrix-vector multiplication. Denotational Semantics

Let:

**v**<sub>car</sub> ∈ 
$$\mathbb{R}^d$$
**M**<sub>fast</sub> ∈  $\mathbb{R}^{d \times d}$ 

Then:

$$[\!["fast car"]\!] = \mathbf{M}_{\mathsf{fast}} \cdot \mathbf{v}_{\mathsf{car}}.$$

### Baroni and Zamparelli: Adjectives as Matrices

#### Type-Theoretic Interpretation

▶ Nouns: *e* : ℝ<sup>*d*</sup>

▶ Adjectives:  $e \rightarrow e : \mathbb{R}^d \rightarrow \mathbb{R}^d$  (linear maps)

Composition:

$$\mathsf{Adj}(\mathsf{N}): \mathbb{R}^d \to \mathbb{R}^d$$
$$\llbracket"\mathsf{fast}" \rrbracket(\llbracket"\mathsf{car}" \rrbracket) = \mathbf{M}_{\mathsf{fast}} \cdot \mathbf{v}_{\mathsf{car}}.$$

イロト イポト イヨト イヨト

# HRR: Binding with Circular Convolution

- Words and roles are hypervectors  $(\mathbb{R}^d)$ .
- Binding is performed using circular convolution.

# Denotational Semantics

Let:

▶ 
$$\mathbf{v}_{car}, \mathbf{v}_{fast} \in \mathbb{R}^d$$
  
Then: [["fast car"]] =  $\mathbf{v}_{fast} \circledast \mathbf{v}_{car}$ , where:  
 $(\mathbf{v}_{fast} \circledast \mathbf{v}_{car})_i = \sum_{i=0}^{d-1} v_{fast,i} \cdot v_{car,(i-j) \mod d}$ .

## HRR: Binding with Circular Convolution

#### Type-Theoretic Interpretation

Nouns:  $e : \mathbb{R}^d$ 

• Adjectives:  $e \to e : \mathbb{R}^d \to \mathbb{R}^d$  (convolution operators) Composition:

$$\mathsf{Adj}(\mathsf{N}): \mathbb{R}^d o \mathbb{R}^d$$
 $[" fast" ]([" car" ]) = \mathbf{v}_{\mathsf{fast}} \circledast \mathbf{v}_{\mathsf{car}}$ 

イロン 不同 とうほう 不同 とう