Information-theoretic locality properties of natural language

Richard Futrell

Department of Language Science Department of Computer Science University of California, Irvine

> @rljfutrell rfutrell@uci.edu

Quantitative Syntax 2019 2019-08-26

• Each human language is a solution to the problem of **maximally efficient communication**...

- Each human language is a solution to the problem of maximally efficient communication...
- subject to fixed human information processing constraints.

- Each human language is a solution to the problem of maximally efficient communication...
- subject to fixed human information processing constraints.
- Efficiency Hypothesis: Languages are optimized so that messages we want to express are easy to produce and comprehend accurately. (Zipf, 1949; Hockett, 1960; Slobin, 1973; Givón, 1991, 1992; Hawkins, 1994, 2004, 2014; Christiansen & Chater, 2008; Jaeger & Tily, 2011; Fedzechkina et al., 2012; MacDonald, 2013)

- Efficiency Hypothesis: Languages are optimized so that messages we want to express are easy to produce and comprehend accurately. (Zipf, 1949; Hockett, 1960; Slobin, 1973; Givón, 1991, 1992; Hawkins, 1994, 2004, 2014; Christian and Comprehence & Title 2011; Federabline et al. 2010;
 - 2014; Christiansen & Chater, 2008; Jaeger & Tily, 2011; Fedzechkina et al., 2012; MacDonald, 2013)

- Efficiency Hypothesis: Languages are optimized so that messages we want to express are easy to produce and comprehend accurately. (Zipf, 1949; Hockett, 1960; Slobin, 1973; Givón, 1991, 1992; Hawkins, 1994, 2004, 2014; Christiansen & Chater, 2008; Jaeger & Tily, 2011; Fedzechkina et al., 2012; MacDonald, 2013)
- Mathematical formalization: human languages are solutions to a constrained optimization problem describing communication subject to cognitive constraints.
 - So, what is the objective function that human languages optimize?

• For example, Ferrer i Cancho & Solé (2003) propose that, for a source random variable *M*, natural languages *L* are minima of:

• For example, Ferrer i Cancho & Solé (2003) propose that, for a source random variable *M*, natural languages *L* are minima of:

$$J_M(L) = H[M | L] + \lambda H[L]$$

• For example, Ferrer i Cancho & Solé (2003) propose that, for a source random variable *M*, natural languages *L* are minima of:

$$J_{M}(L) = H[M|L] + \lambda H[L]$$

Ambiguity of meaning of the signal. (Conditional entropy of meaning given signal)

• For example, Ferrer i Cancho & Solé (2003) propose that, for a source random variable *M*, natural languages *L* are minima of:



• For example, Ferrer i Cancho & Solé (2003) propose that, for a source random variable *M*, natural languages *L* are minima of:



 This function is also known as the Deterministic Information Bottleneck (Strouse & Schwab, 2016) and the Infomax Criterion (Bell & Sejnowski, 1995; Friston, 2010).

• For example, Ferrer i Cancho & Solé (2003) propose that, for a source random variable *M*, natural languages *L* are minima of:



 This function is also known as the Deterministic Information Bottleneck (Strouse & Schwab, 2016) and the Infomax Criterion (Bell & Sejnowski, 1995; Friston, 2010).

Key part: effort is quantified using entropy (average surprisal).

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)
 - Bob threw out the trash.



• The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)



• Bob threw the trash out. \downarrow

• The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)



• Bob threw the trash out.

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)
 - Bob threw out the trash.
 - Bob threw the trash out.
 - Bob threw out the old trash that had been sitting in the kitchen. \downarrow

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)
 - Bob threw out the trash.
 - Bob threw the trash out. \downarrow
 - Bob threw out the old trash that had been sitting in the kitchen. \downarrow

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)
 - Bob threw out the trash.
 - Bob threw the trash out. \downarrow
 - Bob threw out the old trash that had been sitting in the kitchen. \downarrow
 - Bob threw the old trash that had been sitting in the kitchen out. \P

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality** (aka dependency length minimization, dependency distance minimization, domain minimization, early immediate constituents, principle of head proximity, Behaghel's First Law, ...)
 - Bob threw out the trash.
 - Bob threw the trash out. \downarrow
 - Bob threw out the old trash that had been sitting in the kitchen. \downarrow
 - Bob threw the old trash that had been sitting in the kitchen out. \checkmark

• The most powerful efficiency-based model of **word order** in natural language is **dependency locality:**

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality:**
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality:**
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)
 - So the **linear distance** between words in dependencies **should be minimized** (for recent reviews, see Dyer, 2017; Temperley & Gildea, 2018; Liu et al., 2018).

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality:**
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)
 - So the **linear distance** between words in dependencies **should be minimized** (for recent reviews, see Dyer, 2017; Temperley & Gildea, 2018; Liu et al., 2018).
 - Explains pervasive word order patterns across languages:

- The most powerful efficiency-based model of word order in natural language is dependency locality:
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)
 - So the **linear distance** between words in dependencies **should be minimized** (for recent reviews, see Dyer, 2017; Temperley & Gildea, 2018; Liu et al., 2018).
 - Explains pervasive word order patterns across languages:
 - Harmonic word order correlations (Greenberg, 1963; Hawkins, 1994)

- The most powerful efficiency-based model of **word order** in natural language is **dependency locality:**
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)
 - So the **linear distance** between words in dependencies **should be minimized** (for recent reviews, see Dyer, 2017; Temperley & Gildea, 2018; Liu et al., 2018).

• Explains pervasive word order patterns across languages:

- Harmonic word order correlations (Greenberg, 1963; Hawkins, 1994)
- Short-before-long and long-before-short preferences (Hawkins, 1994, 2004, 2014; Wasow, 2002)

- The most powerful efficiency-based model of word order in natural language is dependency locality:
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)
 - So the **linear distance** between words in dependencies **should be minimized** (for recent reviews, see Dyer, 2017; Temperley & Gildea, 2018; Liu et al., 2018).

• Explains pervasive word order patterns across languages:

- Harmonic word order correlations (Greenberg, 1963; Hawkins, 1994)
- Short-before-long and long-before-short preferences (Hawkins, 1994, 2004, 2014; Wasow, 2002)
- Tendency to projectivity (Ferrer-i-Cancho, 2006)

- The most powerful efficiency-based model of word order in natural language is dependency locality:
 - Robust evidence from psycholinguistics that long dependencies cause processing difficulty (Gibson, 1998, 2000; Grodner & Gibson, 2005; Bartek et al., 2011)
 - So the **linear distance** between words in dependencies **should be minimized** (for recent reviews, see Dyer, 2017; Temperley & Gildea, 2018; Liu et al., 2018).

• Explains pervasive word order patterns across languages:

- Harmonic word order correlations (Greenberg, 1963; Hawkins, 1994)
- Short-before-long and long-before-short preferences (Hawkins, 1994, 2004, 2014; Wasow, 2002)
- Tendency to projectivity (Ferrer-i-Cancho, 2006)

Focus of this Work

Focus of this Work

• **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:
 - Information locality: Words are under pressure to be close in proportion to their mutual information.

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:
 - Information locality: Words are under pressure to be close in proportion to their mutual information.
- I show that information locality makes correct predictions beyond dependency locality in two domains:

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:
 - Information locality: Words are under pressure to be close in proportion to their mutual information.
- I show that information locality makes correct predictions beyond dependency locality in two domains:
 - (1) Differences between different dependencies

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:
 - Information locality: Words are under pressure to be close in proportion to their mutual information.
- I show that information locality makes correct predictions beyond dependency locality in two domains:
 - (1) Differences between different dependencies
 - (2) Relative order of adjectives

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:
 - Information locality: Words are under pressure to be close in proportion to their mutual information.
- I show that information locality makes correct predictions beyond dependency locality in two domains:
 - (1) Differences between different dependencies
 - (2) Relative order of adjectives

- **Problem**. Dependency locality is motivated in terms of heuristic arguments about memory usage.
- <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?
- <u>Answer</u>. When we adopt a more sophisticated model of processing difficulty, we can derive dependency locality as a special case of a new information-theoretic principle:
 - Information locality: Words are under pressure to be close in proportion to their mutual information.
- I show that information locality makes correct predictions beyond dependency locality in two domains:
 - (1) Differences between different dependencies
 - (2) Relative order of adjectives

Information Locality

- Introduction
- Information Locality
- Study 1: Strength of Dependencies
- Study 2: Adjective Order
- Conclusion

• **Surprisal theory** (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):

- **Surprisal theory** (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:

- **Surprisal theory** (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)

- Surprisal theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)



- Surprisal theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)
- Accounts for:



- Surprisal theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)
- Accounts for:
 - Garden path effects (Hale, 2001)



- Surprisal theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)
- Accounts for:
 - Garden path effects (Hale, 2001)
 - Antilocality effects (Konieczny, 2000; Levy, 2008)



- Surprisal theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)
- Accounts for:
 - Garden path effects (Hale, 2001)
 - Antilocality effects (Konieczny, 2000; Levy, 2008)
 - Syntactic construction frequency effects (Levy, 2008)



- Surprisal theory (Hale, 2001; Levy, 2008; Smith & Levy, 2013; Hale, 2016):
- Processing difficulty at a word is equal to the surprisal of that word in context:
- Difficulty(w | context) = -logP(w | context)
- Accounts for:
 - Garden path effects (Hale, 2001)
 - Antilocality effects (Konieczny, 2000; Levy, 2008)
 - Syntactic construction frequency effects (Levy, 2008)
- In other words, the average processing difficulty in a language is proportional to the entropy of the language *H*[*L*].



 Surprisal theory has excellent empirical coverage for observable processing difficulty, *except*

- Surprisal theory has excellent empirical coverage for observable processing difficulty, *except*
- It **does not account for dependency locality effects** empirically (Levy, 2008, 2013) and provably cannot theoretically (Levy, 2006; Futrell, 2017).

- Surprisal theory has excellent empirical coverage for observable processing difficulty, *except*
- It **does not account for dependency locality effects** empirically (Levy, 2008, 2013) and provably cannot theoretically (Levy, 2006; Futrell, 2017).
- Reason: Surprisal theory has no notion of **memory limitations**.

- Surprisal theory has excellent empirical coverage for observable processing difficulty, *except*
- It **does not account for dependency locality effects** empirically (Levy, 2008, 2013) and provably cannot theoretically (Levy, 2006; Futrell, 2017).
- Reason: Surprisal theory has no notion of **memory limitations**.
- So how can we build memory limitations into surprisal theory?

- Surprisal theory has excellent empirical coverage for observable processing difficulty, *except*
- It **does not account for dependency locality effects** empirically (Levy, 2008, 2013) and provably cannot theoretically (Levy, 2006; Futrell, 2017).
- Reason: Surprisal theory has no notion of **memory limitations**.
- So how can we build memory limitations into surprisal theory?

- Surprisal theory has excellent empirical coverage for observable processing difficulty, *except*
- It **does not account for dependency locality effects** empirically (Levy, 2008, 2013) and provably cannot theoretically (Levy, 2006; Futrell, 2017).
- Reason: Surprisal theory has no notion of **memory limitations**.
- So how can we build memory limitations into surprisal theory?

Futrell & Levy (2017)

Surprisal: Diff(w | context) = -logP(w | context)

Futrell & Levy (2017)

Surprisal: Diff(w | context) = -logP(w | context)

context	W
Bob threw the old trash that had been sitting in the kitchen	out

Surprisal: Diff(w | context) = -logP(w | context)

context	W
Bob threw the old trash that had been sitting in the kitchen	out



Futrell & Levy (2017)

• Surprisal: $Diff(w | context) = -\log P(w | context)$





Futrell & Levy (2017)

Surprisal: Diff(w | context) = -logP(w | context)





Surprisal: Diff(w | context) = -logP(w | context)


Surprisal: Diff(w | context) = -logP(w | context)





Surprisal: Diff(w | context) = -logP(w | context)





Surprisal: Diff(w | context) = -logP(w | context)





Lossy-context surprisal: Diff(w | context) = -logP(w | memory representation)





Lossy-context surprisal: Processing difficulty per word is

•

Lossy-context surprisal: Processing difficulty per word is

•

$$\mathsf{Diff}(w_i | w_{1,\ldots,i-1}) \propto -\log p(w_i | m_i),$$

Lossy-context surprisal: Processing difficulty per word is

•

$$\mathsf{Diff}(w_i | w_{1,\ldots,i-1}) \propto -\log p(w_i | m_i),$$

where *m_i* is a lossy compression of the context *w*_{1,...,i-1}, i.e. *m_i* is an approximate epsilon-machine (Feldman & Crutchfield, 1998; Marzen & Crutchfield, 2017).

Lossy-context surprisal: Processing difficulty per word is

•

$$\mathsf{Diff}(w_i | w_{1,\ldots,i-1}) \propto -\log p(w_i | m_i),$$

where *m_i* is a lossy compression of the context *w*_{1,...,i-1}, i.e. *m_i* is an approximate epsilon-machine (Feldman & Crutchfield, 1998; Marzen & Crutchfield, 2017).

So the average processing difficulty for a language is a cross entropy:

Lossy-context surprisal: Processing difficulty per word is

•

$$\mathsf{Diff}(w_i | w_{1,\ldots,i-1}) \propto -\log p(w_i | m_i),$$

where *m_i* is a lossy compression of the context *w*_{1,...,i-1}, i.e. *m_i* is an approximate epsilon-machine (Feldman & Crutchfield, 1998; Marzen & Crutchfield, 2017).

So the average processing difficulty for a language is a cross entropy:

$$\mathsf{Diff}(L) \propto \mathbb{E}\left[-\log p(w_i \mid m_i)\right]_{w_{1,\dots,i}}$$

Lossy-context surprisal: Processing difficulty per word is

•

$$\mathsf{Diff}(w_i | w_{1,\ldots,i-1}) \propto -\log p(w_i | m_i),$$

where *m_i* is a lossy compression of the context *w*_{1,...,i-1}, i.e. *m_i* is an approximate epsilon-machine (Feldman & Crutchfield, 1998; Marzen & Crutchfield, 2017).

So the average processing difficulty for a language is a cross entropy:

$$\mathsf{Diff}(L) \propto \mathbb{E} \left[-\log p(w_i \mid m_i)\right]$$
$$\equiv H_I[L']$$

memory representation

Bob threw the old trash sitting in the kitchen	out
--	-----

memory representation

Bob threw the old trash sitting in the kitchen	out
--	-----

memory representation

Bob threw the old trash sitting in the kitchen	out
--	-----

 If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.

memory representation

Bob threw the old trash sitting in the kitchen

out

- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.

memory representation

Bob threw the old trash sitting in the kitchen

out

- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- <u>Theorem</u> (Futrell & Levy, 2017):

memory representation

Bob threw the old trash sitting in the kitchen

- out
- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

memory representation

Bob threw the old trash sitting in the kitchen

- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

: 1

out

memory representation

Bob threw the old trash sitting in the kitchen

- out
- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

Pointwise mutual information (pmi) is the most general statistical measure of *how strongly two values predict each other* (Church & Hanks, 1990) $pmi(w; w') = \log p(w|w')$ p(w)

memory representation

Bob threw the old trash sitting in the kitchen

- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

: 1

out

memory representation

Bob threw the old trash sitting in the kitchen

- out
- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

memory representation

Bob threw the old trash sitting in the kitchen

- out
- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

memory representation

Bob threw the old trash sitting in the kitchen

- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

 e_d : Proportion of information retained about the *d*'th most recent word (Under the noisy memory model, this must decrease monotonically.)

out

memory representation

Bob threw the old trash sitting in the kitchen

- If information about words is lost at a constant rate (noisy memory), then the memory representation will have less information about words that have been in memory longer.
- This leads to information locality. Difficulty increases when words with high mutual information are distant.
- Theorem (Futrell & Levy, 2017):

$$Diff(w_i|w_1, ..., w_{i-1}) \approx -\log P(w) - \sum_{j=1}^{i-1} e_{i-j} pmi(w_i; w_j)$$

 e_d : Proportion of information retained about the *d*'th most recent word (Under the noisy memory model, this must decrease monotonically.)

out

• Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.

- Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.
- How does this relate to **dependency locality**?

- Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.
- How does this relate to **dependency locality**?
- Linking Hypothesis: Words in syntactic dependencies have high mutual information (de Paiva Alves, 1996; Yuret, 1998)

- Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.
- How does this relate to **dependency locality**?
- Linking Hypothesis: Words in syntactic dependencies have high mutual information (de Paiva Alves, 1996; Yuret, 1998)
 - Makes sense a priori: Mutual information is a measure of **strength of covariance**.

- Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.
- How does this relate to **dependency locality**?
- Linking Hypothesis: Words in syntactic dependencies have high mutual information (de Paiva Alves, 1996; Yuret, 1998)
 - Makes sense a priori: Mutual information is a measure of **strength of covariance**.
 - If this is true, then we can see **dependency locality effects as a subset of information locality effects**.

- Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.
- How does this relate to **dependency locality**?
- Linking Hypothesis: Words in syntactic dependencies have high mutual information (de Paiva Alves, 1996; Yuret, 1998)
 - Makes sense a priori: Mutual information is a measure of **strength of covariance**.
 - If this is true, then we can see dependency locality effects as a subset of information locality effects.
- I have a talk about this tomorrow! (Futrell, Qian, Gibson, Fedorenko & Blank, 2019).

- Information locality: I predict processing difficulty when words that predict each other (have high mutual information) are far apart.
- How does this relate to **dependency locality**?
- Linking Hypothesis: Words in syntactic dependencies have high mutual information (de Paiva Alves, 1996; Yuret, 1998)
 - Makes sense a priori: Mutual information is a measure of strength of covariance.
 - If this is true, then we can see **dependency locality effects as a subset of information locality effects**.
- I have a talk about this tomorrow! (Futrell, Qian, Gibson, Fedorenko & Blank, 2019).



- Introduction
- Information Locality
- Study 1: Strength of Dependencies
- Study 2: Adjective Order
- Conclusion

Strength of Dependencies

Strength of Dependencies

• **Dependency locality** says: All words in dependencies should be close.

Strength of Dependencies

- **Dependency locality** says: All words in dependencies should be close.
- Information locality says: Words want to be close in proportion to their mutual information.
- **Dependency locality** says: All words in dependencies should be close.
- Information locality says: Words want to be close in proportion to their mutual information.
- Information locality prediction: Words in dependencies which predict each other other strongly will be especially attracted to each other, beyond dependency locality effects.

- **Dependency locality** says: All words in dependencies should be close.
- Information locality says: Words want to be close in proportion to their mutual information.
- Information locality prediction: Words in dependencies which predict each other other strongly will be especially attracted to each other, beyond dependency locality effects.

$$y_r = \beta_0 + \beta_{\rm pmi} {\rm pmi}(h; d)$$

$$y_r = \beta_0 + \beta_{\rm pmi} {\rm pmi}(h;d)$$

Distance between words in the *r'th* dependency in language *I*



 So: Fit a regression predicting the distance between a head and dependent from the pmi of the head and dependent.



• Fit to UD v2.1 corpora of 50 languages.



- Fit to UD v2.1 corpora of 50 languages.
- I measure pmi between POS tags, not wordforms, because wordform mutual information is hard to estimate for natural language (see my talk tomorrow)

 I find a significant pmi attraction effect in 48/50 languages.

 I find a significant pmi attraction effect in 48/50 languages.

Language	β_{MI}	p	Language	β_{MI}	p	Language	β_{MI}	p
Ancient Greek	-0.18	<.001	Hindi	-0.26	<.001	Slovak	-0.30	<.001
Arabic	-0.26	<.001	Hungarian	-0.11	<.001	Slovenian	-0.38	<.001
Basque	-0.22	<.001	Indonesian	-0.22	<.001	Spanish	-0.37	<.001
Belarusian	-0.20	<.001	Irish	-0.37	<.001	Swedish	-0.35	<.001
Bulgarian	-0.29	<.001	Italian	-0.35	<.001	Tamil	-0.18	<.001
Catalan	-0.29	<.001	Japanese	-0.32	<.001	Turkish	-0.22	<.001
Church Slavonic	-0.23	<.001	Kazakh	-1.18	0.01	Ukrainian	-0.29	<.001
Coptic	-0.35	<.001	Korean	-0.14	<.001	Urdu	-0.22	<.001
Croatian	-0.32	<.001	Latin	-0.18	<.001	Uyghur	-0.04	0.79
Czech	-0.27	<.001	Latvian	-0.32	<.001	Vietnamese	-0.27	<.001
Danish	-0.38	<.001	Lithuanian	-0.41	<.001			
Dutch	-0.10	<.001	Mandarin	-0.19	<.001			
English	-0.38	<.001	Modern Greek	-0.25	<.001			
Estonian	-0.32	<.001	Norwegian	-0.37	<.001			
Finnish	-0.29	<.001	Persian	-0.19	<.001			
French	-0.33	<.001	Polish	-0.23	<.001			
Galician	-0.35	<.001	Portuguese	-0.23	<.001			
German	-0.25	<.001	Romanian	-0.36	<.001			
Gothic	-0.19	<.001	Russian	-0.18	<.001			
Hebrew	-0.21	<.001	Sanskrit	0.10	0.28			

- I find a significant pmi attraction effect in 48/50 languages.
- Average effect size is -0.3:

Language	β_{MI}	p	Language	β_{MI}	p	Language	β_{MI}	p
Ancient Greek	-0.18	<.001	Hindi	-0.26	<.001	Slovak	-0.30	<.001
Arabic	-0.26	<.001	Hungarian	-0.11	<.001	Slovenian	-0.38	<.001
Basque	-0.22	<.001	Indonesian	-0.22	<.001	Spanish	-0.37	<.001
Belarusian	-0.20	<.001	Irish	-0.37	<.001	Swedish	-0.35	<.001
Bulgarian	-0.29	<.001	Italian	-0.35	<.001	Tamil	-0.18	<.001
Catalan	-0.29	<.001	Japanese	-0.32	<.001	Turkish	-0.22	<.001
Church Slavonic	-0.23	<.001	Kazakh	-1.18	0.01	Ukrainian	-0.29	<.001
Coptic	-0.35	<.001	Korean	-0.14	<.001	Urdu	-0.22	<.001
Croatian	-0.32	<.001	Latin	-0.18	<.001	Uyghur	-0.04	0.79
Czech	-0.27	<.001	Latvian	-0.32	<.001	Vietnamese	-0.27	<.001
Danish	-0.38	<.001	Lithuanian	-0.41	<.001			
Dutch	-0.10	<.001	Mandarin	-0.19	<.001			
English	-0.38	<.001	Modern Greek	-0.25	<.001			
Estonian	-0.32	<.001	Norwegian	-0.37	<.001			
Finnish	-0.29	<.001	Persian	-0.19	<.001			
French	-0.33	<.001	Polish	-0.23	<.001			
Galician	-0.35	<.001	Portuguese	-0.23	<.001			
German	-0.25	<.001	Romanian	-0.36	<.001			
Gothic	-0.19	<.001	Russian	-0.18	<.001			
Hebrew	-0.21	<.001	Sanskrit	0.10	0.28			

- I find a significant pmi attraction effect in 48/50 languages.
- Average effect size is -0.3:
 - For each bit of pmi between two words, they are 0.3 words closer together on average.

	Language	$\beta_{\rm MI}$	p	Language	β_{MI}	p	Language	$\beta_{\rm MI}$	p
	Ancient Greek	-0.18	<.001	Hindi	-0.26	<.001	Slovak	-0.30	<.001
	Arabic	-0.26	<.001	Hungarian	-0.11	<.001	Slovenian	-0.38	<.001
	Basque	-0.22	<.001	Indonesian	-0.22	<.001	Spanish	-0.37	<.001
	Belarusian	-0.20	<.001	Irish	-0.37	<.001	Swedish	-0.35	<.001
	Bulgarian	-0.29	<.001	Italian	-0.35	<.001	Tamil	-0.18	<.001
	Catalan	-0.29	<.001	Japanese	-0.32	<.001	Turkish	-0.22	<.001
	Church Slavonic	-0.23	<.001	Kazakh	-1.18	0.01	Ukrainian	-0.29	<.001
	Coptic	-0.35	<.001	Korean	-0.14	<.001	Urdu	-0.22	<.001
	Croatian	-0.32	<.001	Latin	-0.18	<.001	Uyghur	-0.04	0.79
	Czech	-0.27	<.001	Latvian	-0.32	<.001	Vietnamese	-0.27	<.001
	Danish	-0.38	<.001	Lithuanian	-0.41	<.001			
	Dutch	-0.10	<.001	Mandarin	-0.19	<.001			
	English	-0.38	<.001	Modern Greek	-0.25	<.001			
	Estonian	-0.32	<.001	Norwegian	-0.37	<.001			
	Finnish	-0.29	<.001	Persian	-0.19	<.001			
	French	-0.33	<.001	Polish	-0.23	<.001			
	Galician	-0.35	<.001	Portuguese	-0.23	<.001			
	German	-0.25	<.001	Romanian	-0.36	<.001			
	Gothic	-0.19	<.001	Russian	-0.18	<.001			
	Hebrew	-0.21	<.001	Sanskrit	0.10	0.28			

Information Locality

- Introduction
- Information Locality
- Study 1: Strength of Dependencies
- Study 2: Adjective Order
- Conclusion

The pretty red Italian car 👍

The pretty red Italian car 👍

The red pretty Italian car 👎

The pretty red Italian car 👍 The red pretty Italian car 👎 The Italian pretty red car 👎

The pretty red Italian car 🧼 The red pretty Italian car 👎 The Italian pretty red car 👎 The pretty Italian red car 👎

The pretty red Italian car 🤞 The red pretty Italian car 👎 The Italian pretty red car 👎 The pretty Italian red car 👎

The pretty red Italian car 🦆 The red pretty Italian car 👎 The Italian pretty red car 👎 The pretty Italian red car 👎







 There are constraints on relative order of adjectives that are stable across speakers and languages.



- There are constraints on relative order of adjectives that are stable across speakers and languages.
- Strongest empirical generalization: more *subjective* adjectives are farther out (Scontras et al., 2017)



- There are constraints on relative order of adjectives that are stable across speakers and languages.
- Strongest empirical generalization: more *subjective* adjectives are farther out (Scontras et al., 2017)
- Information locality explanation: Adjectives with high pmi with a noun will appear relatively close to that noun.



- There are constraints on relative order of adjectives that are stable across speakers and languages.
- Strongest empirical generalization: more *subjective* adjectives are farther out (Scontras et al., 2017)
- Information locality explanation: Adjectives with high pmi with a noun will appear relatively close to that noun.
 - Possibly conceptually related to subjectivity.

• 1. Gather a large set of adjective—adjective—noun triples from a corpus.

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?
- Data: Google Syntactic n-Grams (8.5 billion adjective-noun pairs)

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?
- Data: Google Syntactic n-Grams (8.5 billion adjective-noun pairs)
- Model: Logistic regression predicting order from pmi.

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?
- Data: Google Syntactic n-Grams (8.5 billion adjective-noun pairs)
- Model: Logistic regression predicting order from pmi.
- Result: PMI predicts adjective order for held-out data with 66.9% accuracy.
Does Adjective Order Correspond to Mutual Information?

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?
- Data: Google Syntactic n-Grams (8.5 billion adjective-noun pairs)
- Model: Logistic regression predicting order from pmi.
- Result: PMI predicts adjective order for held-out data with 66.9% accuracy.
 - Best previously known predictor (subjectivity) gets 68.4%

Does Adjective Order Correspond to Mutual Information?

- 1. Gather a large set of adjective—adjective—noun triples from a corpus.
- 2. Measure pmi between adjectives and nouns.
- 3. Does pmi(A;N) predict that A will be closer to the noun than the other adjective?
- Data: Google Syntactic n-Grams (8.5 billion adjective-noun pairs)
- Model: Logistic regression predicting order from pmi.
- **Result: PMI predicts adjective order** for held-out data with **66.9%** accuracy.
 - Best previously known predictor (subjectivity) gets 68.4%
 - PMI + Subjectivity gets 72.9% accuracy

• Other theories aim to explain the same data...

- Other theories aim to explain the same data...
 - Dyer's (2017, 2018) **Integration Cost**: Involves the conditional entropy of dependency relation labels given words.

- Other theories aim to explain the same data...
 - Dyer's (2017, 2018) Integration Cost: Involves the conditional entropy of dependency relation labels given words.
 - Hahn et al.'s (2018) Subjective Rational Speech Acts Model: Involves noisy incremental memory in the computation of meaning.

- Other theories aim to explain the same data...
 - Dyer's (2017, 2018) Integration Cost: Involves the conditional entropy of dependency relation labels given words.
 - Hahn et al.'s (2018) Subjective Rational Speech Acts Model: Involves noisy incremental memory in the computation of meaning.
 - Scontras et al.'s (2019) **Noisy composition model** explains adjective order in terms of noisy hierarchical computation of meaning.

- Other theories aim to explain the same data...
 - Dyer's (2017, 2018) Integration Cost: Involves the conditional entropy of dependency relation labels given words.
 - Hahn et al.'s (2018) Subjective Rational Speech Acts Model: Involves noisy incremental memory in the computation of meaning.
 - Scontras et al.'s (2019) **Noisy composition model** explains adjective order in terms of noisy hierarchical computation of meaning.
- Future work will have to rigorously disentangle the predictions of these theories.

- Other theories aim to explain the same data...
 - Dyer's (2017, 2018) Integration Cost: Involves the conditional entropy of dependency relation labels given words.
 - Hahn et al.'s (2018) Subjective Rational Speech Acts Model: Involves noisy incremental memory in the computation of meaning.
 - Scontras et al.'s (2019) **Noisy composition model** explains adjective order in terms of noisy hierarchical computation of meaning.
- Future work will have to rigorously disentangle the predictions of these theories.
- Problem: The relevant information-theoretic quantities are hard to estimate accurately.

Information Locality

- Introduction
- Information Locality
- Study 1: Strength of Dependencies
- Study 2: Adjective Order
- Conclusion

• <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?

• <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?

$$J_M(L) = H[M | L] + \lambda H[L]$$

• <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?

$$J_M(L) = H[M | L] + \lambda H_L[L']$$

 <u>Question</u>. How does dependency locality fit in formally with information-theoretic models of natural language?

 $J_M(L) = H[M | L] + \lambda H_L[L']$ **Dependency** locality happens in this term in the form of

information locality

 Does information locality capture the trade-off of complex morphology and deterministic word order? (Koplenig et al., 2017)

- Does information locality capture the trade-off of complex morphology and deterministic word order? (Koplenig et al., 2017)
 - Depends on the precise relationship between morphology and inter-word MI.

- Does information locality capture the trade-off of complex morphology and deterministic word order? (Koplenig et al., 2017)
 - Depends on the precise relationship between morphology and inter-word MI.
- Is the right notion of mutual information purely MI between words, or is it also something that takes into account meaning?

- Does information locality capture the trade-off of complex morphology and deterministic word order? (Koplenig et al., 2017)
 - Depends on the precise relationship between morphology and inter-word MI.
- Is the right notion of mutual information purely MI between words, or is it also something that takes into account meaning?
 - E.g., dependency relation types, as in Dyer's Integration Cost theory

- Does information locality capture the trade-off of complex morphology and deterministic word order? (Koplenig et al., 2017)
 - Depends on the precise relationship between morphology and inter-word MI.
- Is the right notion of mutual information purely MI between words, or is it also something that takes into account meaning?
 - E.g., dependency relation types, as in Dyer's Integration Cost theory
- Does information locality make different predictions from dependency locality wrt crossing dependencies?

- Does information locality capture the trade-off of complex morphology and deterministic word order? (Koplenig et al., 2017)
 - Depends on the precise relationship between morphology and inter-word MI.
- Is the right notion of mutual information purely MI between words, or is it also something that takes into account meaning?
 - E.g., dependency relation types, as in Dyer's Integration Cost theory
- Does information locality make different predictions from dependency locality wrt crossing dependencies?

Thanks all!

- All code is available online at <u>http://github.com/langprocgroup/adjorder</u> and <u>http://github.com/langprocgroup/cliqs</u>
- Thanks to Roger Levy, Ted Gibson, and Tim O'Donnell for discussions.
- Thanks to the SyntaxFest reviewers for helpful comments.
- Thanks to the Quasy organizers for a great conference!