


The Computability of Language Structure

(人類語言中的數學性質)

PeterWolf
Droidtown Linguistic Tech.

Who am I?

- Full-time job: Droidtown Linguistic Tech. founder & general manager.
 - Webpage: <https://api.droidtown.co>
 -  Discord: <https://discord.gg/g5Enb5zAyK>
- Inventor of linguistic technologies.
 - 3 Taiwan invention patents, and 1 US invention patent holder.
- Part-time lecturer in NYCU:
 - Linguistics and Artificial Intelligence.
 - Linguistic Technology and Complex System
 - https://github.com/PeterWolf-tw/LxTech_and_ComplexSys/



We will be talking about
MATH
today.



When it comes to "Math in languages"...

NCCU Corpus of Spoken Taiwan Mandarin

政治大學中文口語語料庫

Home About The Corpus Corpus Data

HOME / About the corpus

About the corpus

The NCCU Corpus of Spoken Taiwan Mandarin is a language documentation whereby open access to the corpus has been collecting spoken data from daily face-to-face conversations obtained from the participants for the publication of an English letter. A broad transcription of speech and code-switching. The spoken data may change over time.

Part of the corpus data are also available at Talkbank: <http://ca.talkbank.org/access/TaiwanMandarin.html>

Fundings for this language documentation project include:

- The Aim for the Top University and Elite Research Center
- The Humanities Research Center of the National Central University
- The Office of Research and Development, National Central University
- Research projects, the Ministry of Science and Technology

NCCU Corpus of Spoken Taiwan Mandarin

政治大學中文口語語料庫

Home About The Corpus Corpus Data Statistics Citations Contact

HOME / Word frequency

> Character frequency

> Word frequency

27 conversational excerpts total about 600 minutes of talk.

2016 詞頻統計 (Frequency count of words)

序次	詞	出現次數	佔所有字元之百分比	累計百分比
1	就	3548	4.261%	4.261%
2	我	3305	3.969%	8.231%
3	啊	3000	3.603%	11.834%
4	是	2792	3.353%	15.187%
5	的	2092	2.513%	17.699%
6	你	1938	2.328%	20.027%
7	那	1685	2.024%	22.051%
8	他	1540	1.850%	23.900%
9	然後	1458	1.751%	25.651%
10	有	1359	1.632%	27.283%

<https://spokentaiwanmandarin.nccu.edu.tw/word-frequency.html>

[網站介紹](#) [網站總覽](#) [語言學習區](#) [語言教學資源區](#) [語料庫資源](#) [研究團隊](#) [相關連結](#) [意見回饋](#) [成果發表](#) [首頁](#) [English](#)

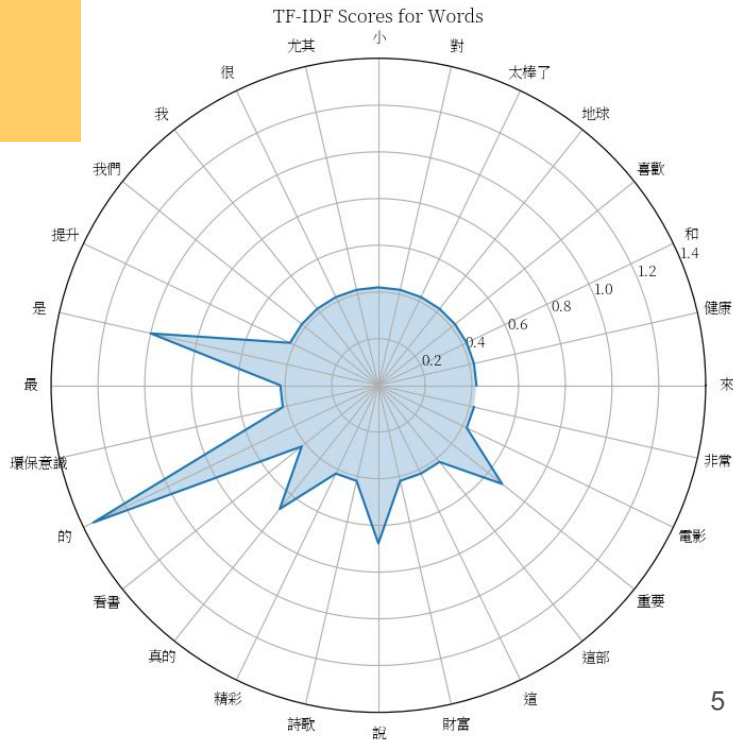


現代漢語語料庫詞頻統計

現代漢語語料庫詞頻統計提供平衡語料庫的詞頻信息。華文教師可依據詞頻統計提供的訊息得知詞語的數量與頻率，從而決定詞語學習的先後安排，幫助教師們編寫教程。

詞頻: TF

表示詞在文檔中出現的頻率，就統計學而言，只要這個詞在文本中出現越多次代表越值得關注，因此它會具有一個重要的統計評估指標之一，但並不是完全相信此統計方式，看完底下的IDF就會知道為什麼。



國家級警報

現在



[防空警報]中國於15:04發射衛星，已飛越南部上空，請民眾注意安全。若發現不明物體，通報警消人員處理。[Air raid Alert] Missile flyover Taiwan airspace, be aware. 國防部 (MND) 02-27355979

<https://vocus.cc/article/659e8a65fd897800012fefbf>

可是為什麼很多人看成「飛越越南」？

從語言學的角度來看，一個很重要的因素發揮了影響力——**詞頻**。

首先，「飛越南部」可以拆成三個詞：

- 飛越
- 越南
- 南部

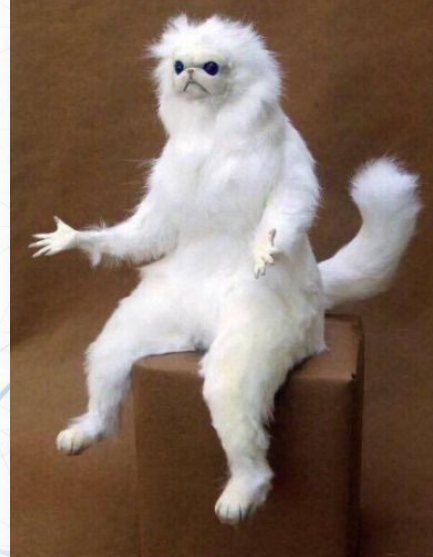
23.21 >>> 3.66，詞頻越高，人類可以越快理解這個詞，所以在「飛越」跟「越南」之間，就先取得了「越南」的資訊。

既然我們是由左至右讀，為什麼不會先抓到「飛越」這個資訊，然後順勢往下看到「南部」，組合成正確的意思？

根據中研院平衡語料庫，這三個詞的詞頻如下：

- 飛越：3.66 (每百萬字會看到 3.66 這個詞)
- 越南：23.21

Frequency!
What a
lovely
word!



Final term paper:

請依以下兩張文字雲，從人、事、時、地、物...等角度寫出 A4 兩頁的分析報告。



背景一：這兩張文字雲是高雄市長選舉結束後，由候選人發表的演講稿製成
(但可能是不同屆哦！)



背景二：這兩個講稿，一個是勝選人講的，一個是敗選人講的



[illegible]

文字雲？腦補 + 暗示的「文本分析」

其實剛剛的文字雲是從這兩篇我自己亂寫的東西生出來的。

所以剛才都是你自己受引導以後的腦補！

一-277萬一個一天一座一樣不夠不好意思之前了人口今天他鄉任
但是偉大充滿內再次到剛剛勇敢包容向哪裡困難國瑜團結報告好
如子弟孕育對就工作已經市民朋友希望幫幫助很很多很大忘恭喜
意思愛跟感受應該打拼拜託挑戰提到時間更好替朋友期永遠熱情
特別當市長疼惜相信禁結束給綠繼續美好而已能致電藍表示許多
誠請變得變故鄉責任走起跟著路辜負這這塊這座這段選舉還有都
長大陪伴陳其邁雖然難過韓市長願望高雄市他他作為作為其邁其
邁努力努力四年四年家園家園從從或或所有所有故鄉故鄉明天明
天更更最後最後生命生命相依相依這裡這裡開始開始一起一起
一起也也也以以以共同共同共同加油加油加油地方地方地方市民
市民市民為榮為榮為榮管管管能夠能夠能夠要要要還是還是還是城
市城市城市城市感謝感謝感謝感謝未來未來未來未來為為為為人
人人人人市長市長市長市長市長跟跟跟跟跟在在在在在支持支
持支持支持支持支持支持不不不不不不不我們我們我們我們我
們我們我們我們大家大家大家大家大家大家大家大家是是是
是是是是是是我我我我我我我我我我我高雄高雄高雄高雄高雄
高雄高雄高雄高雄高雄高雄高雄高雄高雄高雄高雄的的的的的
的的的的的的的的的的的的的的的

—2018年—280萬一件一個一個一個一個月一定一定一定一幕一瓶一絲一起
一起一起一起一遍一點上下一心上陣下一代下來不不不不不停不用不
用不用世界世界中中中央也也事人人人人民今天份佛教徒你來保佑做做
做做價值先內內內內人全全全全力全力以赴全部兩年再出來出現分分
分族群到前前前剛剛才力量動員包容千斤可台灣台灣各位各位向喜
悅因為國民黨團結團隊在在在在城市報告報告報告外多少大家大
家大家大家大家太多太多好不好好朋友如果姊妹委員
委員宣言將對手對手對手小內閣就是是工作已經市民市市民市市民
民市市長市長市長希望幫忙府廉潔廣大建立影往往往很清楚後後得心必
須志工朋友愛感謝感謝感謝感謝感謝感謝感謝感謝感謝感謝成長我
我我我我我我我我我我我我我我我我我我我們我們我們我們
我們我們所以所有所有所有所有所有所有有打拚拼打造批判批判過
找到找到把投票拚拚挑戰挽起提名擔心擔心支持改變政府政治效率敗選
教育敬新明天是晚安更好最好最困難最神奇最重要最高會會會有朝
有陰朋友們未來未來未來未來機會機會民主政治史氣水沉重沒有
的的的的的的的的的的的的的的的的的的的的的的的的的的的的的
盡速直接相信相信相信相信盾盾看看到真眼礦泉祝福票數秉持立刻等待組織
給給經濟經濟老兄弟肩扛肩膀能自己與菩薩華人萬斤藍綠藍綠藍綠表情
表情袖子要得要得要親愛訊號訊請謝謝謝謝謝謝謝謝謝讓讓走
走走起來跟跟跟跟跳出來輕裝逐步這個這次這裡造勢過程選舉選舉
選舉選舉那都鄉親鄉親鄉親鄉親朋友釋放重新重組重複開始關係
關心院長陳其邁陳其邁陳其邁集合非常非常巨大非常強烈鞠躬願意高
市府高雄高雄高雄高雄高雄高雄高雄高雄高雄高雄高雄高雄高雄
高雄高雄高雄高雄高雄市高雄市高雄市高雄市點黨部

可是為什麼很多人看成「飛越越南」？

從語言學的角度來看，一個很重要的因素發揮了影響力——**詞頻**。

首先，「飛越南部」可以拆成三個詞：

- 飛越
- 越南
- 南部

既然我們是由左至右讀，為什麼不會先抓到「飛越」這個資訊，然後順勢抓到「南部」，組合成正確的意思？

根據中研院平衡語料庫，這三個詞的詞頻如下：

- 飛越：3.66 (每百萬字會看到 3.66 這個詞)
- 越南：23.21

但為什麼不是先取得最高詞頻「南部」的資訊，然後正確理解這段話呢？

先問一個問題，「越南」跟「南部」，你覺得哪個詞的所指較精確、較具體、較容易想像？

我的答案是「越南」，如果同意的話可以繼續往下看。

語意學中有主體 (figure) 跟背景 (ground) 的概念，主體較有形、範圍較精確，背景則相反。眼睛比較容易快速看到主體，並且記住他。

既然「越南」比較具體，那麼就是因為這個原因，我們會比較先注意到詞頻高且看似為這段話主體的「越南」，而不是詞頻超高，但看似為背景的「南部」。

所以**主體-背景**的關係會凌駕於詞頻的影響嗎？答案是，我不知道，我還不知道。

AI 模式

全部

圖片

影片

購物

新聞

書籍

更多 ▾

工具 ▾



ACL Anthology

<https://aclanthology.org> > 2020.emnlp... · 翻譯這個網頁 · <https://aclanthology.org/2020.emnlp-main.331/>

Word Frequency Does Not Predict Grammatical ...

由 C Yu 著作 · 2020 · 被引用 18 次 — We find that across four orders of magnitude, corpus frequency is unrelated to a noun's performance on grammatical tasks.



arXiv

<https://arxiv.org> > cs · 翻譯這個網頁 ·

Statistical patterns of word frequency suggesting the ...

由 S Yu 著作 · 2020 · 被引用 1 次 — Abstract page for arXiv paper 2012.00187: Statistical patterns of word frequency suggesting the probabilistic nature of human languages.



Wikipedia

<https://en.wikipedia.org> > Zipf's law · 翻譯這個網頁 ·

Zipf's law

In many texts in human languages, word frequencies approximately follow a Zipf distribution with exponent. At the low-frequency end, where the rank approaches ...



Airiti Library 華藝線上圖書館

<https://www.airitilibrary.com> > Publication > Index :

使用文字探勘實作新聞事件追蹤= News event tracking using ...

2025年8月18日 -- 本論文利用R 語言建立一個新聞事件追蹤系統，透過網路爬蟲爬取新聞文章，將爬取的文章做清理，利用jieba 斷詞後，依據各文章中斷詞的結果建立詞頻矩陣，透過TF-IDF 的計算找出 ...



佛教慈濟醫療財團法人

<https://dlweb01.tzuchi.com.tw> > eclass > CJFD PDF :

中國期刊全文資料庫中國博士學位論文全文數據庫中國優秀 ...

2025年7月13日 -- 詞頻：檢索詞在相應檢索項中出現的 頻次。詞頻為空，表示至少出現1次，. 如果為數字，例如2，則表示至少出 現2次，以此類推. 精確：檢索結果完全等同或包 含與檢索字 ...



中華民國圖書館學會

<https://colisp.lac.org.tw> > pdf > CoLISP2024 PDF :

圖書資訊學術與實務研討會會議論文集

2024年12月14日 -- ... 資訊技術運用於人文學研究上，產生跨領域合作. 的「數位人文」 (digital humanity) 研究。學者藉 助數位工具可針對大量的文本內容進行詞頻計算，. 或是挖掘與探索某些 ...
265 頁



國立政治大學心理學系

<https://psy.nccu.edu.tw> > PageStaffing > Detail :

蔡介立副教授 - 國立政治大學心理學系

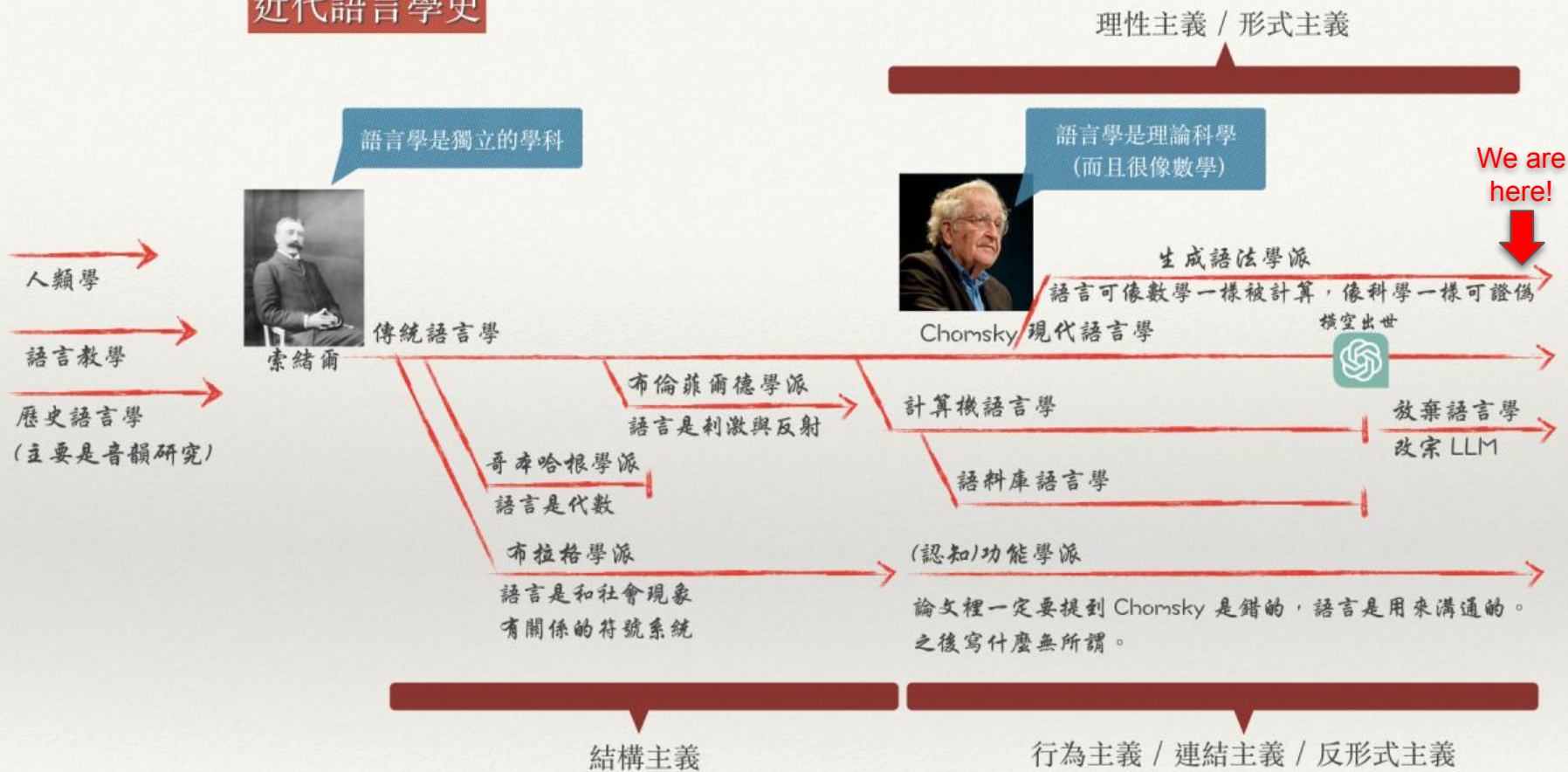
2025年6月15日 -- 第二語言經驗對閱讀第一和第二語言詞頻與詞預測力效果之調控：眼動同步記錄的大腦造影研究(延攬); 蔡介立; 計畫主持人; 2020年08月~2021年07月; 國科會. 年度; 計畫名稱 ...

So, when it comes to frequency...

誰跟你說第一招先用「詞頻」
誰就是在欺負你不懂語言學！



近代語言學史



If you are creating a language...

A language must have "**morphemes**" as its basic units to form words.

A language must have "**syntax**" to govern how words are put together.

A language must have "**semantics**" to convey meaning.

A language must have "**performance**" to deliver information (inward or outward).

Demo: Language Creator

Assuming that we have a language...part 1

- A language must have "**morphemes**" as its basic units to form words.
 - It only has 10 morphemes.

Numeral	
-	0
↖	1
(2
↙	3
⋈	4
ℱ	5
⋈	6
↘	7
⋈	8
↑	9

form words

0123 => 123

246 => 246

855 => 855

56183 => 56183

Quiz:

Is **00000** a possible word in this language?

Assuming that we have a language...part 2

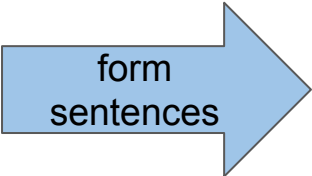
- A language must have "**syntax**" to govern how words are put together.
 - It has two "functional words": +, -
 - It has embedding clauses.

0123 => 123

246 => 246

855 => 855

56183 => 56183



form
sentences

123 + 8550 - 24

46 + 246 - 55

85 + 8 - 5

5 + (6 - 1 + (8 + 3))

Fun facts:

In some dialects, "A + B" is noted as "A B +",
in other dialects, "A + B" is noted as "+ B A."

Quiz: Is this grammatical in this language?

5 + (6 - 1 + (8) + 3)

Assuming that we have a language... part 3

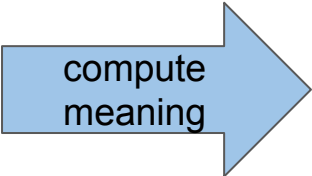
- A language must have "**semantics**" to compute the meaning when words are put together.
 - It has two "functional words": +, -, x
 - It has embedding clauses.

$$123 + 8550 - 24$$

$$46 + 246 - 55$$

$$85 + 8 - 5$$

$$5 + (6 - 1 + (8 + 3))$$



compute
meaning

$$123 + 8550 - 24 = 8649$$

$$46 + 246 - 55 = 237$$

$$85 + 8 - 5 = 88$$

$$5 + (6 - 1 \times (8 + 3)) = 0$$

Quiz: Do you think we come up with the number after "=" by some computational process or some probability model?

Assuming that we have a language... part 4

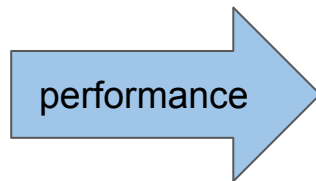
- A language must have "**performance**" to deliver information (inward or outward).

$$123 + 8550 - 24 = 8649$$

$$46 + 246 - 55 = 237$$

$$85 + 8 - 5 = 88$$

$$5 + (6 - 1 \times (8 + 3)) = 0$$



$$123 + 8550 - 24 = 8649$$

-> 8 **thousand** 6 **hundred** and 4**ty**-9

-> 8**ty**-6 4**ty**-9

-> 8 6 4 9

Quiz: Do you agree with the idea that

"MATH is just another LANGUAGE?"

We do have a language...part 1

- A language has "morphemes" as its basic units to form words.

03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	UTF8
𠂇 5903	𠂈 5904	𠂉 5905	𠂊 5906	𠂋 5907	𠂌 5908	𠂍 5909	𠂎 590A	𠂏 590B	𠂐 590C	𠂑 590D	𠂒 590E	𠂓 590F	E5A480
𠂔 5913	𠂕 5914	𠂖 5915	𠂗 5916	𠂘 5917	𠂙 5918	𠂚 5919	𠂛 591A	𠂜 591B	𠂝 591C	𠂞 591D	𠂟 591E	𠂠 591F	E5A490
𠂡 5923	𠂢 5924	𠂣 5925	𠂤 5926	𠂥 5927	𠂦 5928	𠂧 5929	𠂨 592A	𠂩 592B	𠂪 592C	𠂫 592D	𠂬 592E	𠂭 592F	E5A4A0
𠂮 5933	𠂯 5934	𠂰 5935	𠂱 5936	𠂲 5937	𠂳 5938	𠂴 5939	𠂵 593A	𠂶 593B	𠂷 593C	𠂸 593D	𠂹 593E	𠂺 593F	E5A4B0
𠂻 5943	𠂼 5944	𠂽 5945	𠂾 5946	𠂿 5947	𠃀 5948	𠃁 5949	𠃂 594A	𠃃 594B	𠃄 594C	𠃅 594D	𠃆 594E	𠃇 594F	E5A580
𠃈 5953	𠃉 5954	𠃊 5955	𠃋 5956	𠃌 5957	𠃍 5958	𠃎 5959	𠃏 595A	𠃐 595B	𠃑 595C	𠃒 595D	𠃓 595E	𠃔 595F	E5A590
𠃕 5963	𠃖 5964	𠃗 5965	𠃘 5966	𠃙 5967	𠃚 5968	𠃛 5969	𠃜 596A	𠃝 596B	𠃞 596C	𠃟 596D	𠃠 596E	𠃡 596F	E5A5A0
𠃢 5973	𠃣 5974	𠃤 5975	𠃥 5976	𠃦 5977	𠃧 5978	𠃨 5979	𠃩 597A	𠃪 597B	𠃫 597C	𠃬 597D	𠃭 597E	𠃮 597F	E5A5B0
𠃯 5983	𠃰 5984	𠃱 5985	𠃲 5986	𠃳 5987	𠃴 5988	𠃵 5989	𠃶 598A	𠃷 598B	𠃸 598C	𠃹 598D	𠃺 598E	𠃻 598F	E5A680
𠃼 5993	𠃽 5994	𠃾 5995	𠃿 5996	𠄀 5997	𠄁 5998	𠄂 5999	𠄃 599A	𠄄 599B	𠄅 599C	𠄆 599D	𠄇 599E	𠄈 599F	E5A690
𠄉 59A3	𠄊 59A4	𠄋 59A5	𠄌 59A6	𠄍 59A7	𠄎 59A8	𠄏 59A9	𠄐 59AA	𠄑 59AB	𠄒 59AC	𠄓 59AD	𠄔 59AE	𠄕 59AF	E5A6A0
𠄖 59B3	𠄗 59B4	𠄘 59B5	𠄙 59B6	𠄚 59B7	𠄛 59B8	𠄜 59B9	𠄝 59BA	𠄞 59BB	𠄟 59BC	𠄠 59BD	𠄡 59BE	𠄢 59BF	E5A6B0
𠄣 59C3	𠄤 59C4	𠄥 59C5	𠄦 59C6	𠄧 59C7	𠄨 59C8	𠄩 59C9	𠄪 59CA	𠄫 59CB	𠄬 59CC	𠄭 59CD	𠄮 59CE	𠄯 59CF	E5A780

form words

𠂇、𠂈 => ??
 𠂮、𠂯 => 大夫
 𠂒、𠂓 => 夏天
 𠂮、𠂯、𠂯 => 大姊姊

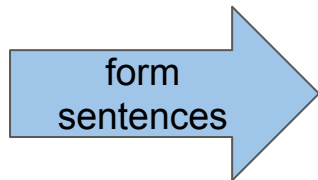
Quiz:

Why "𠂇𠂈" is unlikely to be a possible word in this language?

Assuming that we have a language...part 2

- A language must have "syntax" to govern how words are put together.
 - It has two "functional words": +, -
 - It has embedding clauses.

夯、奄 => ??
大、夫 => 大夫
夏、天 => 夏天
大、姊、姊 => 大姊姊



大姊姊 在 夏天 當上 大夫

大姊姊 聽到 (妹妹 叫 (弟弟 去 游泳))

Fun facts:

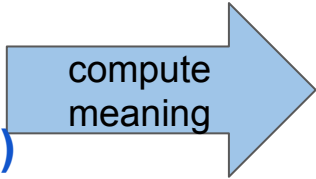
In some languages, "姐姐 當上 大夫" is noted as "姐姐 大夫 當上",
in other languages, "姐姐 當上 大夫" is noted as "當上 大夫 姐姐"

Assuming that we have a language...part 3

- A language must have "semantics" to compute the meaning when words are put together.
 - It has two "functional words": +, -, x
 - It has embedding clauses.

大姊姊 在 夏天 當上 大夫

大姊姊 聽到 (妹妹 叫 (弟弟 去 游泳))



compute
meaning

[[大姊姊在夏天當上大夫]]
= the condition in which
there is a 姊姊 and a 大夫
in between a 當
上-relation holds in the
time when it is 夏天

Assuming that we have this language...part 4

- What are inward/outward "performance"?

[[大姊姊在夏天當上大夫]]

= the condition in which
there is a 姊姊 and a 大夫
in between a 當
上-relation holds in the
time when it is 夏天



performance

You can either "think of" it in
your mind, "speak it out" with
your mouth or sign it with a sign
language.

Fun facts:

Since "speak it out/sign it" are only part of the possible
performance options, communication is not the core part of
language system.

Language = Math

The core of a language system is a set of intertwined computational processes describing...

- How words are formed with morphemes
- How sentences are formed with words?
- How meanings are computed with sentences?
- How the results are performed?
- **Language is a math system with more than the numbers from 0 to 9.**

Why did computational linguistics choose frequency as the tool?

- With the birth of computers in the **1940s**, computational linguistics started.
- Chomsky's first paper came out at late **1950s**.
- When computational linguists started to investigate languages, the only tool in hand was...basically counting words as "**word frequency**."
- Many students nowadays confuse the order of the years and think that "Chomsky is doing traditional linguistics" and "computational linguists are doing modern linguistics." Na...Computational linguistics came first and the approaches are still pretty much the same as their ancestors in the 40s.

Do frequency models reflect human minds?

Large language models are said to be a compression model of the world. ([ref.](#))

- **Language Modeling Is Compression** 
Gregoire Deletang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, Joel Veness
Published: 16 Jan 2024, Last Modified: 15 Mar 2024 ICLR 2024 poster Everyone Revisions BibTeX
Code Of Ethics: I acknowledge that I and all co-authors of this work have read and commit to adhering to the ICLR Code of Ethics.
Keywords: lossless compression, arithmetic coding, language models, scaling laws, in-context learning
Submission Guidelines: I certify that this submission complies with the submission instructions as described on <https://iclr.cc/Conferences/2024/AuthorGuide>.

-  NOTONLY AI
<https://notonly-ai.com> > blog-detail > open-ai-首席技術...
Open AI 首席技術官Ilya Sutskever-GPT是壓縮是壓縮全世界 ...
2023年5月10日 — 神經網絡與GPT"有一種誤解，認為ChatGPT是一個大型語言模型，但有一個圍繞它的系統，"黃仁勳說。Sutskever表示，OpenAI使用兩個級別的訓練 ...

- 
ANNALS OF ARTIFICIAL INTELLIGENCE
CHATGPT IS A BLURRY JPEG OF THE WEB
OpenAI's chatbot offers paraphrases, whereas Google offers quotes. Which do we prefer?
By Ted Chiang
February 9, 2023

Two ways of compression:

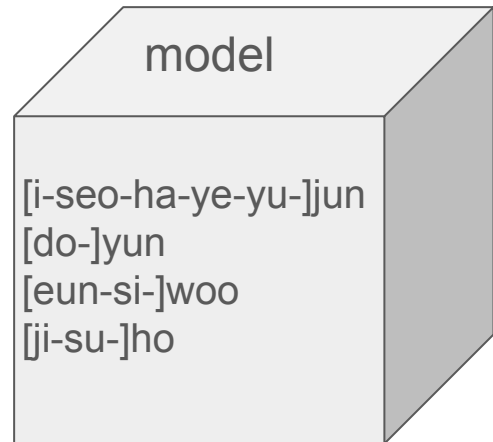
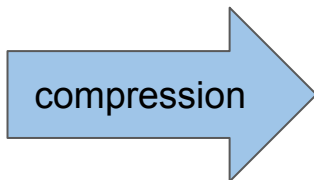
- Frequency-based compression => zip up the "frequently-seen" elements.

List of the most popular given names in South Korea

2021 [\[edit \]](#)

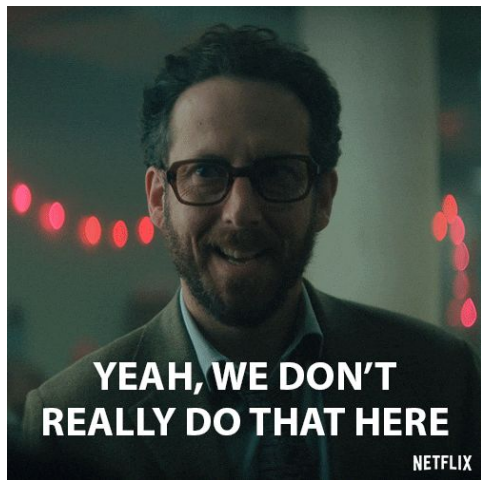
Boys ^[2]				
Common spelling ↕	Hangul ↕	RR ↕	MR ↕	Count ↕
I-jun	이준	Ijun	Ijun	2,833
Seo-jun	서준	Seojun	Sōjun	2,396
Ha-jun	하준	Hajun	Hajun	2,227
Do-yun	도윤	Doyun	Toyun	2,199
Eun-woo	은우	Eunu	Ŭnu	1,931
Si-woo	시우	Siu	Siu	1,831
Ji-ho	지호	Jiho	Chiho	1,606
Ye-jun	예준	Yejun	Yejun	1,455
Yu-jun	유준	Yujun	Yujun	1,380
Su-ho	수호	Suho	Suho	1,360

69 characters in length

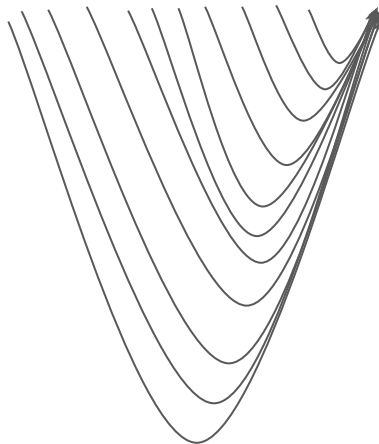


50 characters in length

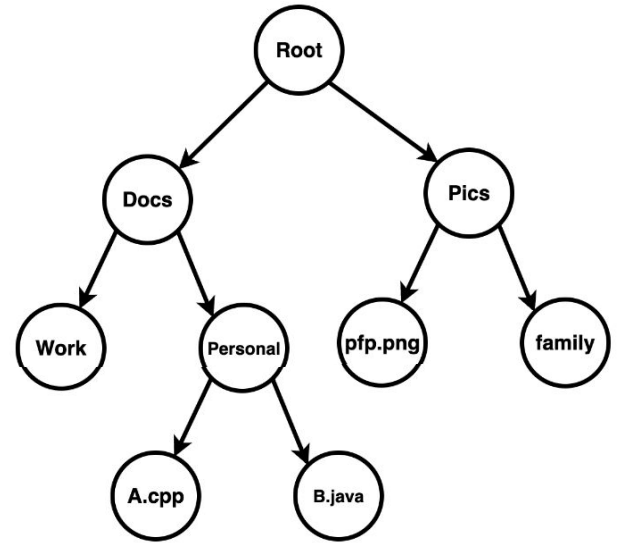
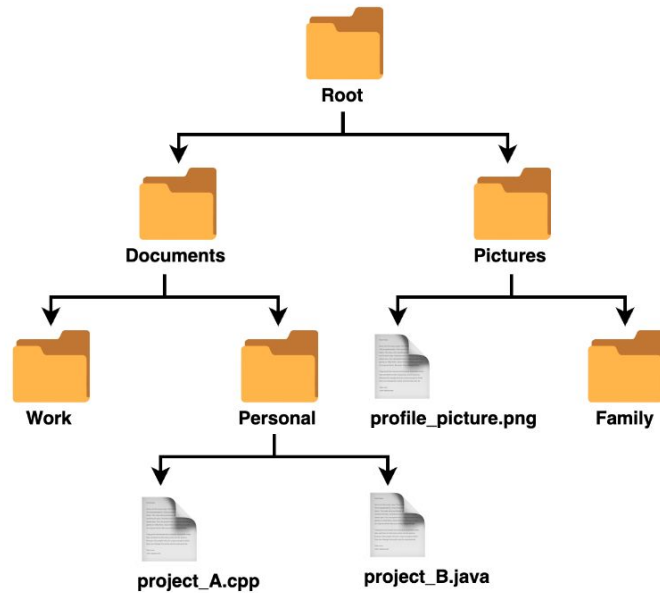
Why frequency model works, and why it is **NOT** how humans do.



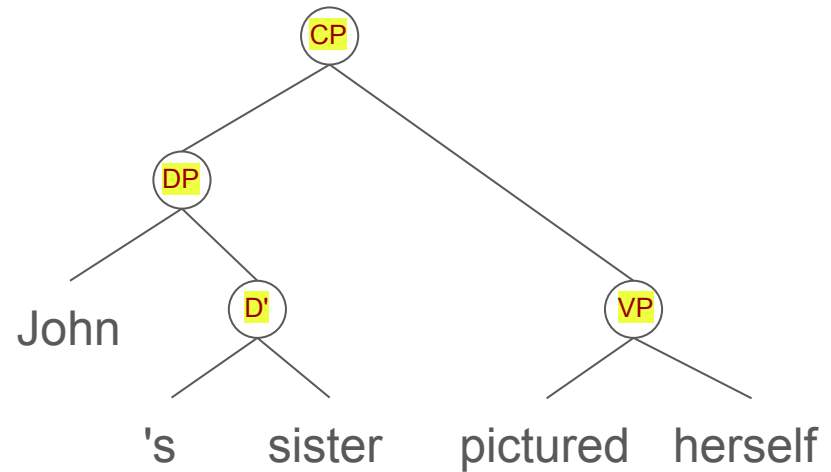
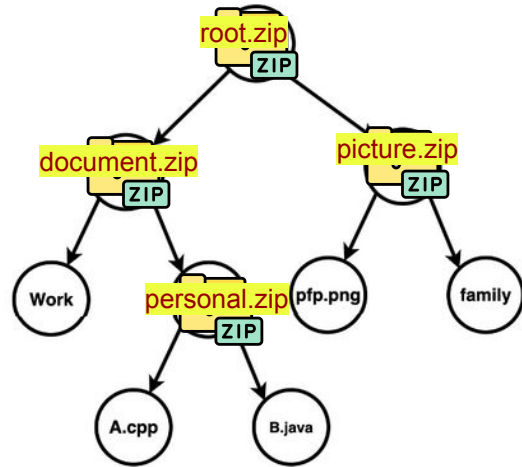
- To be able to "compute" a frequency means you need "a lot of samples."
 - Human children don't have that.
- A frequency model "predicts the next token" **by probability** and that is why LLMs are bad at doing math.
 - $123 + 8550 - 24 = 8649$



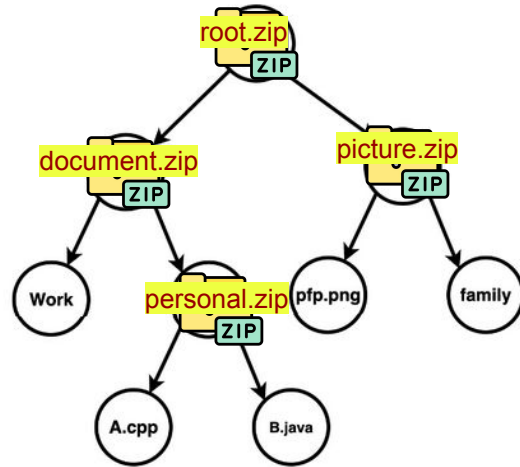
Another kind of compression: Binary-Branching Syntax Tree



Another kind of compression: Binary-Branching Syntax Tree



A compressed object must be decompressible!



 => {document.zip, picture.zip}

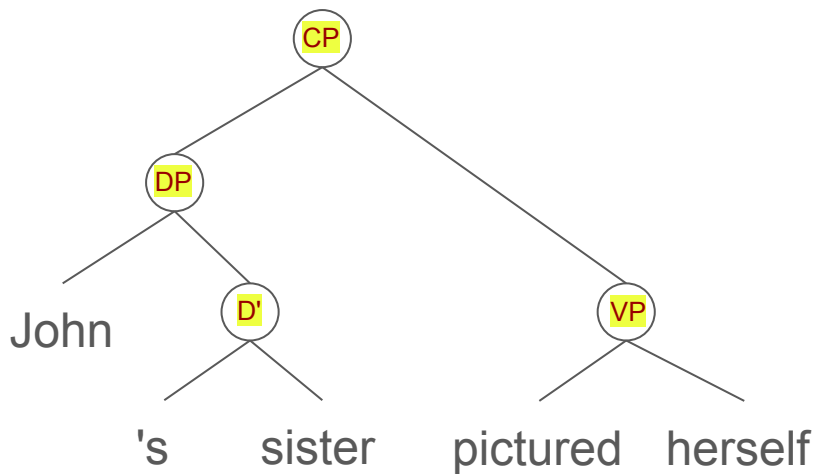
 => {work, personal.zip}

 => {A.cpp, B.java}

 => {pfp.png, family}

- To make sure what are decompressed is the same as what were compressed before, an **ALGORITHM** that can guarantee the output instead of a **probability model** that has no guarantees to the output must be implemented.

A Binary-Branching Syntax Tree is the Algorithm!



CP => {DP*, VP}

DP => {John, D'}

D' => {'s, sister}

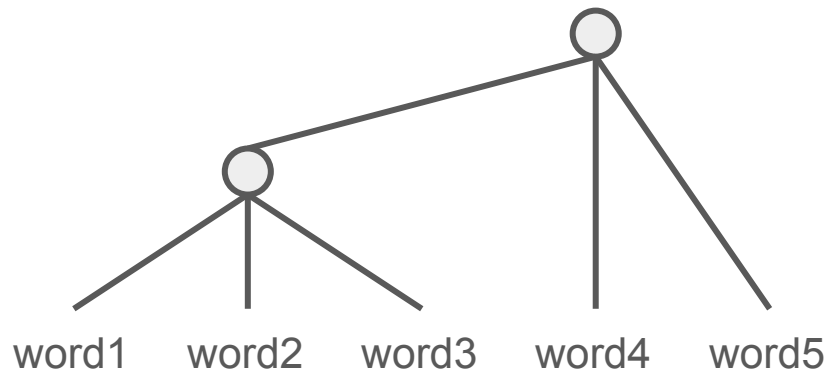
VP => {pictured, herself}

- Don't get misled! This is not the phrase-structure rules from the last century in many aspects:
 - It's **binary**: memory cost at every level is fixed.
 - It's **universal**: the compression/decompression algorithm applies to all human languages.
 - It's **corpus-free**: no data model is needed for training.
 - It's **efficient**: (to be discussed later)
 - one more thing. it is only taught in NYCU in Taiwan.

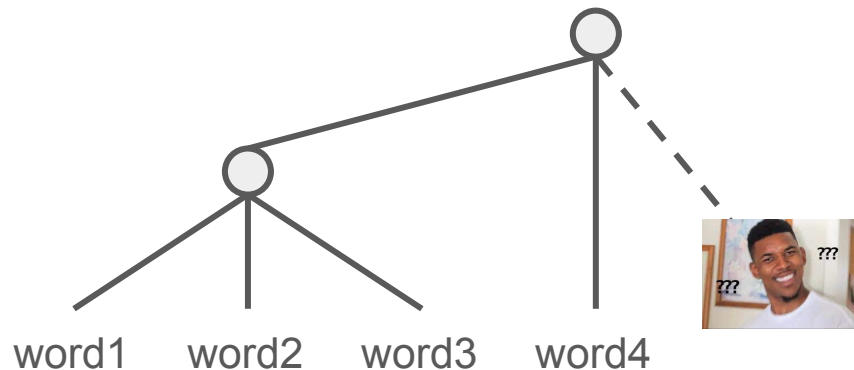
Why not ternary (or more)?

Because you cannot know how many words will be used in a sentence before you perceive it.

- Assuming you have a ternary tree in your head...



Ternary tree works perfect with a 5-word sentence.

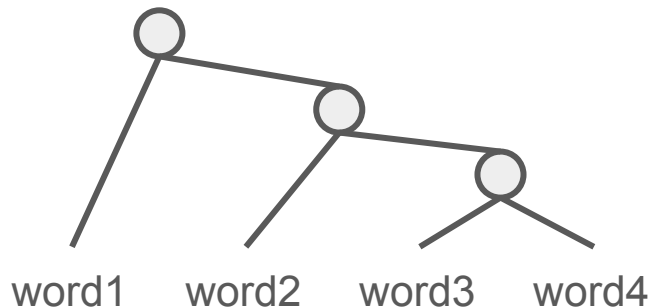
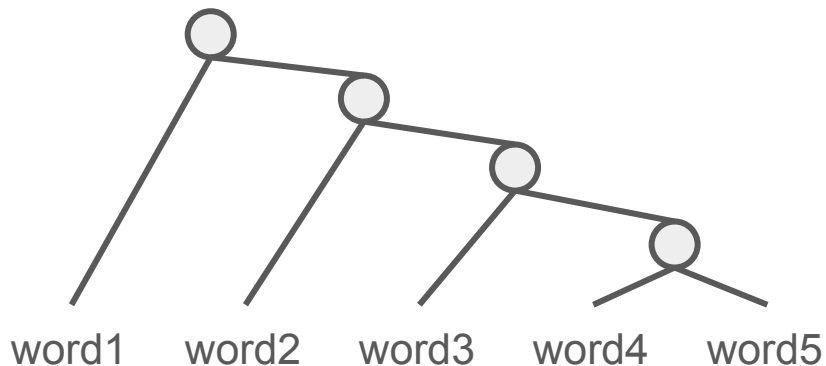


But in a 4-word sentence...

On the other hand, a binary-branching structure...

No matter how many words will be used in a sentence, a binary-branching structure can handle it.

- Assuming you have a ternary tree in your head...



Why binary?

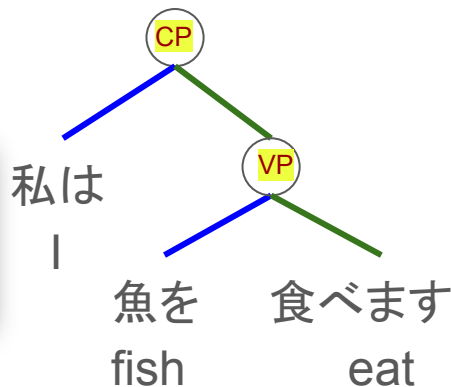
- When designing an algorithm, fixed data length is generally faster than non-fixed data length in computing because it simplifies and speeds up operations like memory access and instruction decoding.
- A fixed-binary-branching tree can handle sentences of different lengths.
- Fixed-ternary (or N-ary where $N > 2$) branching tree can not handle sentences of different lengths.
- Given the fact that human languages use sentences of different length, binary-branching tree is the only viable structure.

Why "Being Universal" is important?

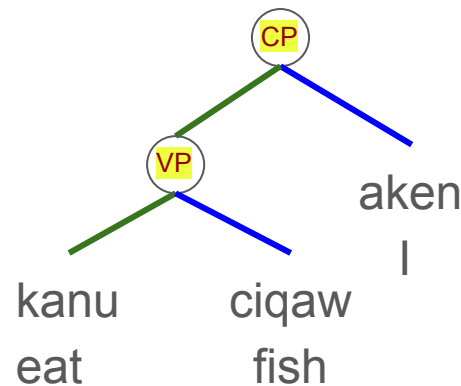
- We are all Homo sapiens; we are the same species.
- Universality guarantees the potentiality of a human to acquire/learn **ANY** human language in the world.

Fun facts:

In some dialects, "A + B" is noted as "A B +",
in other dialects, "A + B" is noted as "+ B A."



Japanese



Paiwan

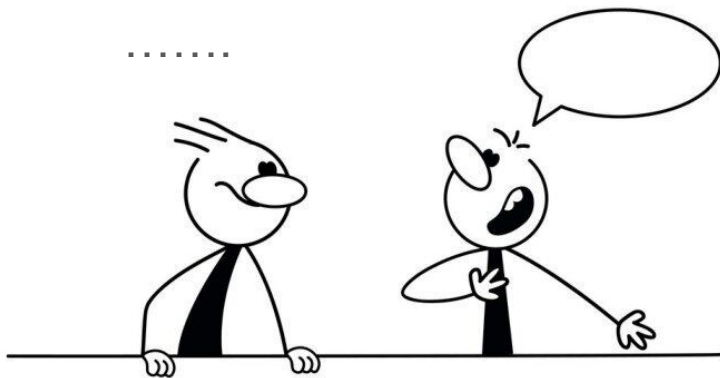
Why do we say binary-branching structure is innate?



When a human starts to use language to communicate with another human being, the first assumption would be that

"WE MUST SHARE THE SAME DECODING/ENCODING SYSTEM."

We don't need to have the same experience in order to have the same language model in mind.



LLM Transfer learning, how promising it is?

License: arXiv.org perpetual non-exclusive license
arXiv:2501.11496v1 [cs.CL] 20 Jan 2025

Generative AI and Large Language Models in
Language Preservation: Opportunities and Challenges

Vincent Koc

Vincent Koc is with Hyperthink, Sydney, Australia (hyperthink.com.au). Contact: vincentkoc@ieee.org.

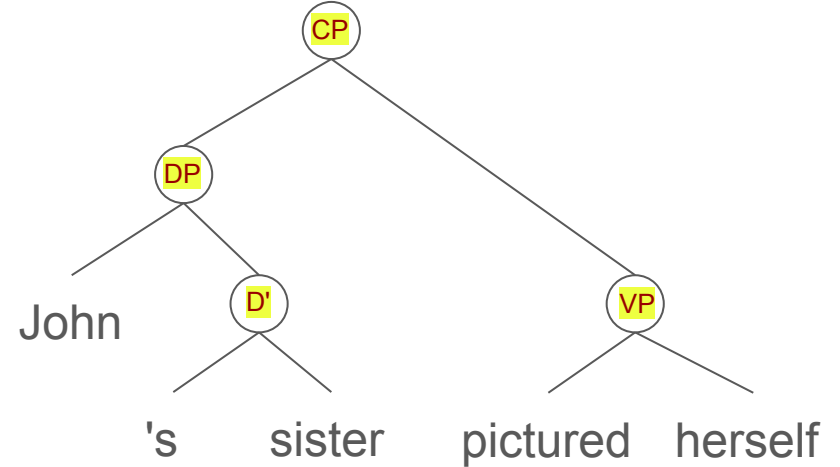
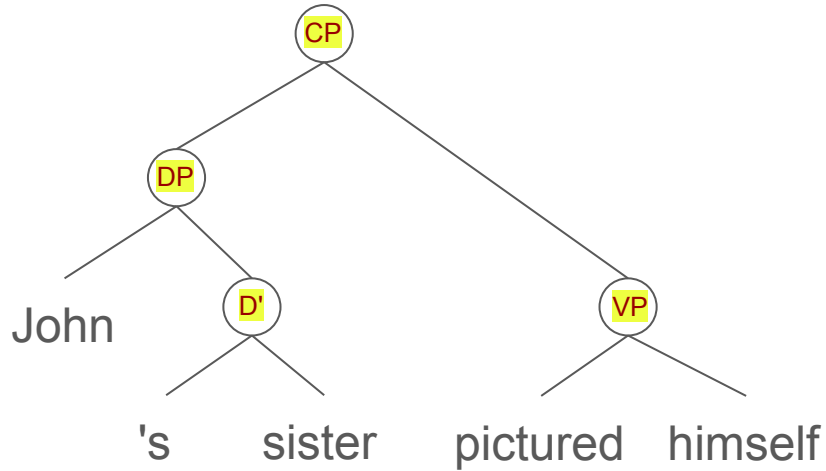
<https://arxiv.org/html/2501.11496v1>

IV-B *Technical Limitations*

Training large-scale language models requires significant computational resources, including specialized hardware and large amounts of electricity. Access to the right infrastructure can be an obstacle, especially for endangered languages spoken in low-resource regions. Additionally, AI models often struggle with complex grammar, non-standard spellings, or extensive lexical borrowing from dominant regional languages, leading to less-accurate outputs.

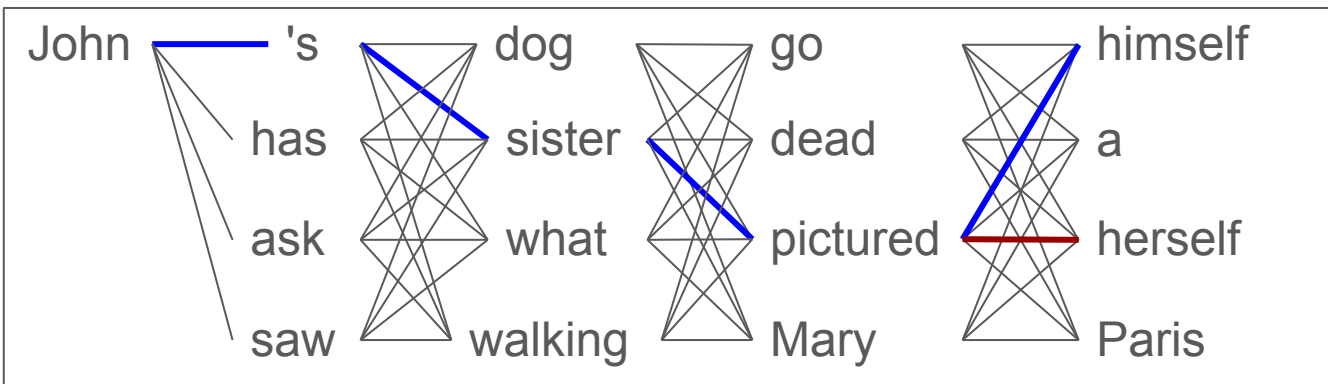
- LLM stands for "LARGE language model", it requires a "LARGE" amount of training data (and electricity) . Endangered languages usually don't have this luxury.
- LLM/AI is not a universal approach.
- If it digitizes/generates less-accurate outputs, then just please don't start. Once wrong corpus is established/generated, it will be there forever.

Efficiency: Which one is grammatical?



Two ways to determine whether "himself" is grammatical here

- Let's train a language model to predict whether it is more frequently to see "himself" or "herself" here.

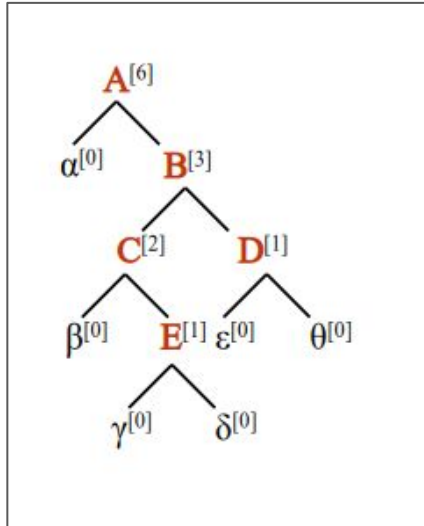


- Why is **"herself"** grammatical but **"himself"** ungrammatical here?
 - Because the probability of "herself" is higher than "himself." (←That's description, not explanation.)
- Why is **"a"** ungrammatical here? Don't you think **"pictured a ... "** would be more frequent than **"pictured herself"** in any "well-balanced" corpus?

Another kind of computation: C-Command

Definition of C-Command:

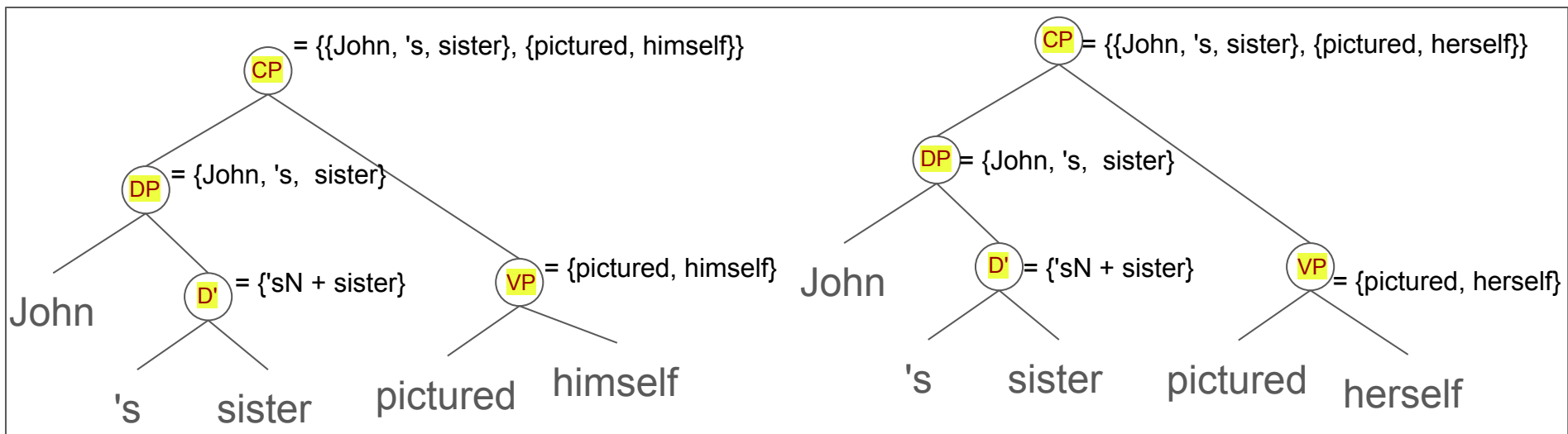
- Node α c-commands a node β iff
- Neither α nor β dominates the other, and
- Every branching node dominating α also dominates β



Red texts are compression nodes.

Computation with Algorithm:

- Why is "**herself**" grammatical but "**himself**" ungrammatical here?
 - Referent of **pronominal expressions** needs to be resolved.
 - The compression node containing "**herself**" is compressed with {John, 's, sister} but not John.
 - The compression node containing "**himself**" is compressed with {John, 's, sister} but not John.



- Is it just some gender-agreement issue? No, it's the structure issue! Consider this:
 - Mary's **sister** took a **selfie** by the tower. (Who is in the selfie? Mary or the sister?)

Comparison:

- Frequency-based Approach

- a. EXTREMELY large training data is required.
- b. MASSIVE electricity is required for computing.
- c. DIFFICULT to take care of endangered languages.
- d. The output is NOT EXPLAINABLE.
- e. Needs a lot of \$\$ investment, like...

A LOT!

- Algorithm-based Approach

- a. Minimum data is needed (at least one instance).
- b. My old 2017 macbook air can do the computation.
- c. See (a). above.
- d. The algorithmic logic is transparent and explainable. Explainable AI means Responsible AI.
- e. Need linguists who know the structures of languages to do the programming.

Takeaways

1. Language as an Algorithmic System

- Language operates like mathematics—through computational processes, not probability

2. Binary-Branching Syntax Trees as Compression

- Human language structure uses **binary-branching trees** as a universal compression/decompression algorithm
- This is the only structure that can handle variable-length sentences efficiently with fixed memory costs
- Universal across all human languages (e.g., Japanese, Paiwan, English all use the same underlying structure)

3. Practical Implications

- The C-command algorithm explains grammaticality (e.g., "herself" vs "himself") through structural relationships, not frequency
- LLMs are not suitable for endangered language preservation—they risk generating inaccurate corpus that persists forever
- True language digitization requires understanding linguistic structure, not just statistical patterns.

Discussions



shutterstock.com • 2556379069