



CS274A - Natural Language Processing



Spring 2024

Administrative Stuff

- ▶ Instructor: Kewei Tu (屠可伟)
 - ▶ Email: tukw@shanghaitech.edu.cn
 - ▶ Office: SIST 1A-304B
 - ▶ Office hours: by appointment
- ▶ TA: 吴昊一、惠文阳、吉鹏宇
 - ▶ Office hours: TBA



Administrative Stuff

▶ Classes

- ▶ Tue/Thu 8:15-9:55am @教学中心201
- ▶ 12 weeks

▶ Prerequisite

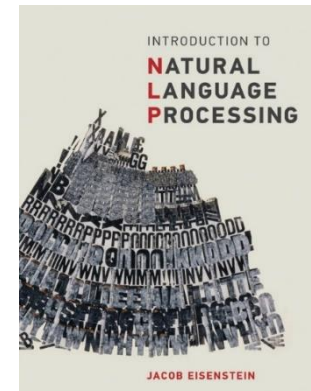
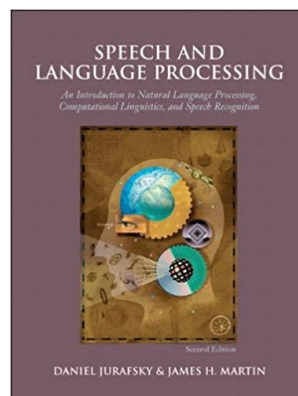
- ▶ CS: Programming, Data Structures and Algorithms
- ▶ Math: Calculus, Probability and Statistics, Linear Algebra
- ▶ Artificial Intelligence I (*recommended*)



Administrative Stuff

► Textbooks

- [SLP] Speech and Language Processing, by Daniel Jurafsky and James Martin
 - 2nd edition published in 2008. 中译版：《自然语言处理综论（第二版）》
 - 3rd edition draft can be found online (updated on Feb 3, 2024)
- [INLP] Introduction to Natural Language Processing, by Jacob Eisenstein

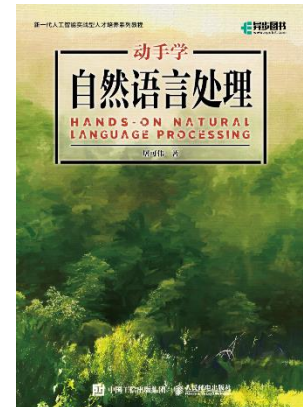


Administrative Stuff

▶ Textbooks

▶ [DSX] 动手学NLP

- ▶ More consistent with this course
- ▶ Contains executable code (Jupyter notebook)
- ▶ To appear this year



Administrative Stuff

- ▶ Blackboard
 - ▶ Announcements, slides, homework assignments, etc.
- ▶ Piazza
 - ▶ Discussion and QA
 - ▶ <https://piazza.com/shanghaitech.edu.cn/spring2024/cs274a/>
- ▶ AutoLab
 - ▶ Project



Administrative Stuff

▶ Grading

- ▶ Homework (20%): 6 homework assignments, due in 7 days
- ▶ Final (60%): late May
- ▶ Project (20%): after week 12



Administrative Stuff

▶ Plagiarism

- ▶ All assignments must be done individually
 - ▶ You may not look at solutions from any other source
 - ▶ You may not share solutions with any other students
 - ▶ Plagiarism detection software will be used on all the programming submissions
- ▶ Way of collaboration
 - ▶ You may discuss together or help another student debug code; however, you cannot dictate or give the exact solution



Administrative Stuff

- ▶ Plagiarism punishment
 - ▶ When one student copies from another student, both students are responsible
 - ▶ Zero point on the assignment
 - ▶ Repeated violation will result in an F grade for this course as well as further discipline at the school/university level





A Brief Introduction to NLP











What is the hottest AI system today?



- ▶ Released on Nov. 30, 2022
- ▶ Huge impact in academia, industry and general public

HOW LONG IT TOOK TOP APPS TO HIT 100M MONTHLY USERS

ChatGPT is estimated to have hit 100M users in January, 2 months after it's launch. Here's how long it took other top apps to reach that:

APP	MONTHS TO REACH 100M GLOBAL MAUS
 CHATGPT	2
 TIKTOK	9
 INSTAGRAM	30
 PINTEREST	41
 SPOTIFY	55
 TELEGRAM	61
 UBER	70
 GOOGLE TRANSLATE	78

SOURCE: UBS

yahoo!
finance

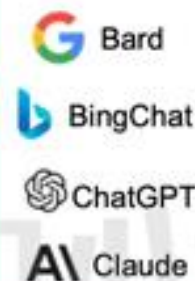
ChatGPT & large language models

国外

基础模型



ChatBot



其他应用



基础模型



ChatBot



其他应用



国内

What is the hottest AI system today?



请写一段自然语言处理研究生课程的开场白



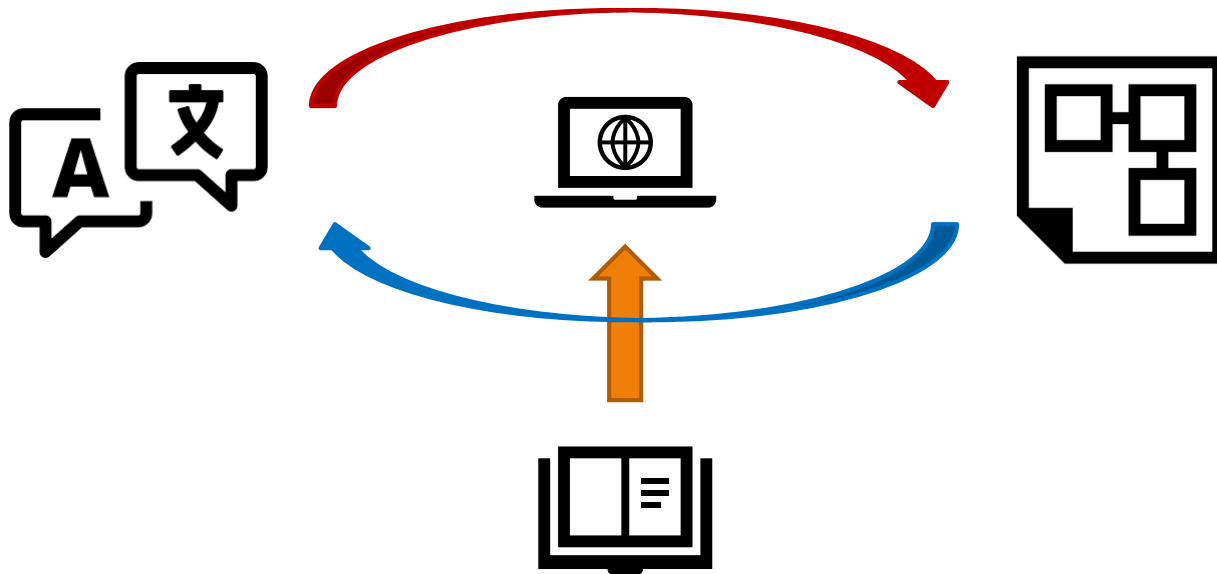
欢迎大家来到自然语言处理研究生课程。自然语言处理是当今科技发展的热门领域之一，在人工智能、大数据和认知科学的背景下发展迅速。本课程将从理论和实践的角度对自然语言处理的基本概念、方法和应用进行深入的讲解。希望各位通过本课程的学习，能够对自然语言处理有更深刻的了解，提高自己的技术水平，同时也能激发各位对这个领域的兴趣。本课程将通过课堂讲解、实验练习、研究项目等多种形式进行教学，希望大家能积极参与，多提问题，多交流，一起探讨自然语言处理的未来。那么，让我们开始这段充满机会与挑战的旅程吧！

*ChatGPT is a (prominent) example of **Natural Language Processing (NLP)**, the subject of this course!*

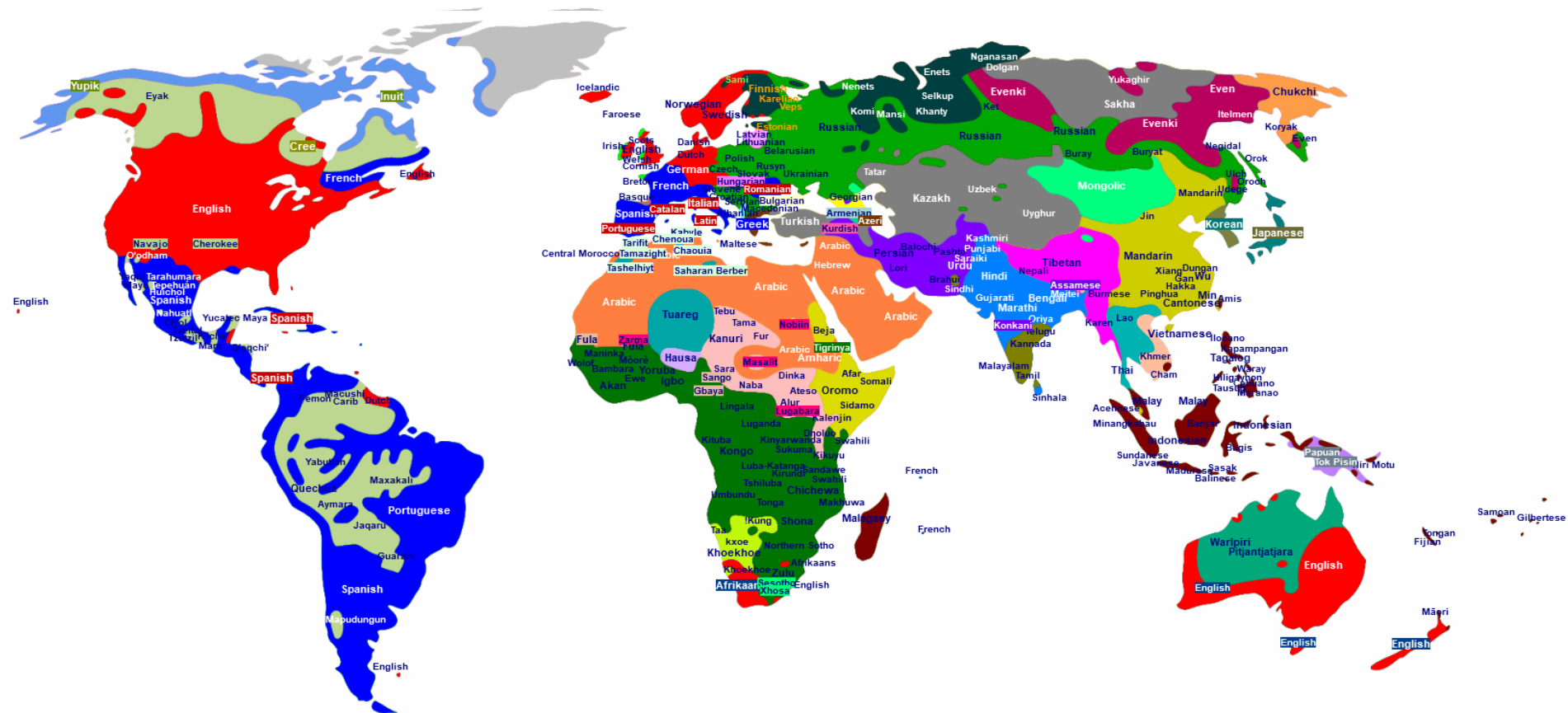
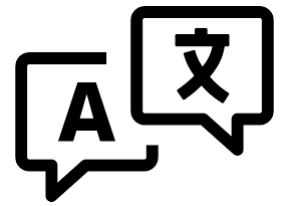


What is NLP?

- ▶ Automating the **analysis**, **generation**, and **acquisition** of human (“natural”) language



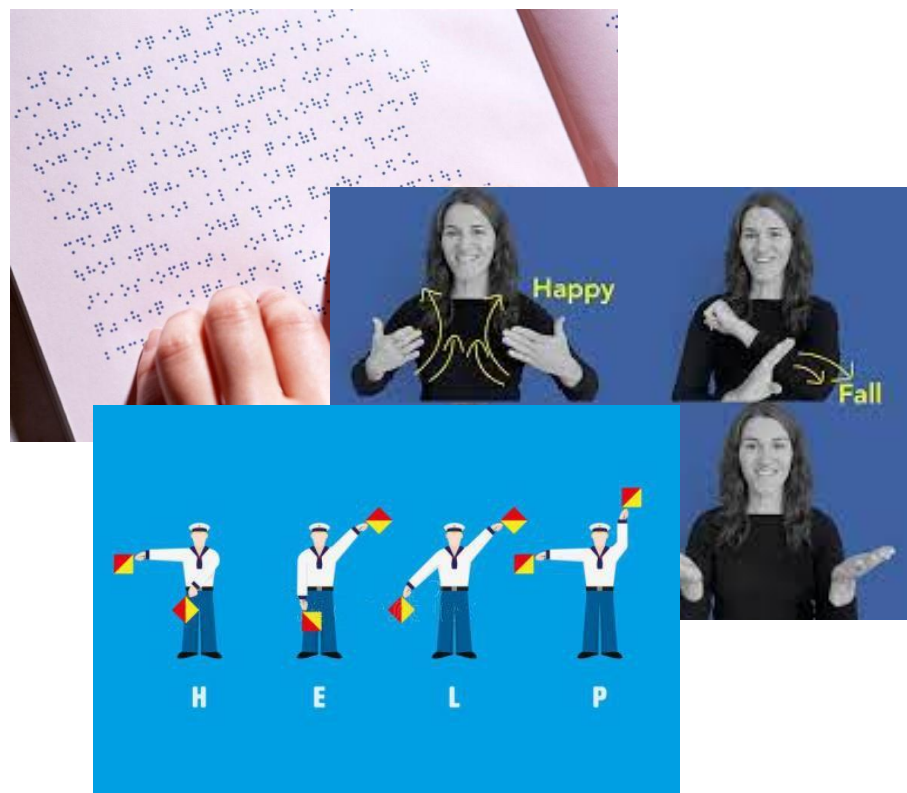
Which language?

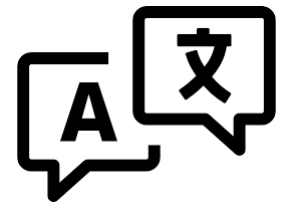




Which language?

- ▶ Ideally, NLP is language-neutral
 - ▶ NLP technology can be applied to any human language
 - ▶ ...if its text can be represented as a sequence of symbols





Which language?

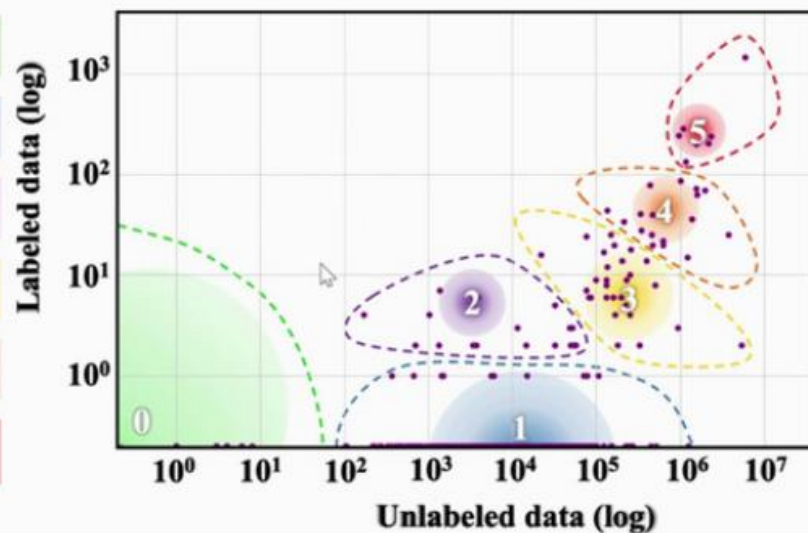
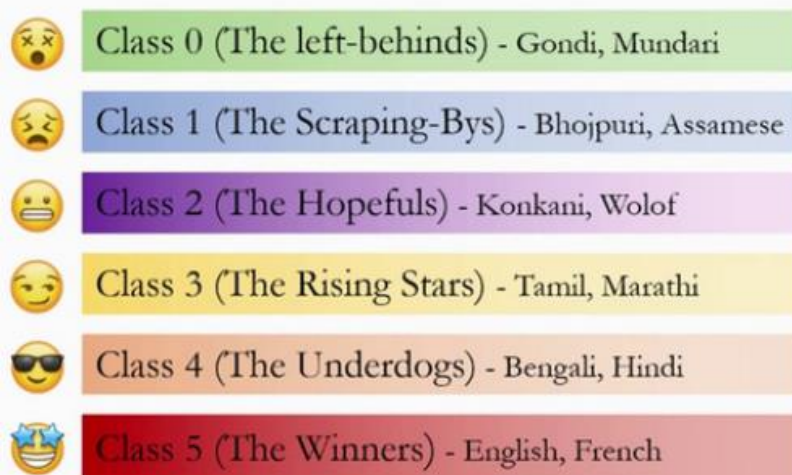
- ▶ Ideally, NLP is language-neutral
 - ▶ NLP technology can be applied to any language
 - ▶ ...if its text can be represented as a sequence of symbols
- ▶ In reality, NLP for some languages is better developed
 - ▶ More interest
 - ▶ Users, market, ...
 - ▶ More resources
 - ▶ Developers, data, computers, \$, ...





Which language?

- ▶ Ideally, NLP is language-neutral
 - ▶ NLP technology can be applied to any language
 - ▶ ...if its text can be represented as a sequence of symbols
- ▶ In reality, NLP for some languages is better developed





What representation?

- ▶ Ideally, a formal language that is sufficiently expressive
 - ▶ First-order predicate logic
 - ▶ Programming language
 - ▶ Neural (distributed) representations??
- ▶ In reality, depends on the application
 - ▶ Labels, features, commands, ...



Fields related to NLP

- ▶ Machine learning
 - ▶ ML is a powerful (but not the only) tool in NLP
 - ▶ NLP is a source of inspiration for ML
- ▶ Linguistics
 - ▶ Roughly: science vs. engineering
 - ▶ NLP \Leftrightarrow computational linguistics
- ▶ Artificial intelligence
 - ▶ NLP is a subfield of AI
 - ▶ “NLP is the crown jewel of AI”
 - ▶ Solving NLP requires solving strong AI



Fields related to NLP

- ▶ Speech Processing
 - ▶ Largely separate from NLP
 - ▶ but there is some overlap
- ▶ Cognitive science / Neuroscience
 - ▶ Humans: the only working NLP prototype!
- ▶ Logic, knowledge representation & reasoning
 - ▶ NLP analyzes NL to and generates NL from logic language
- ▶ Theory of computation
 - ▶ Studies formal language and grammars
 - ▶ Provides a lot of tools to NLP



NLP Applications

► Chatbot

Assistants



Chit-Chat

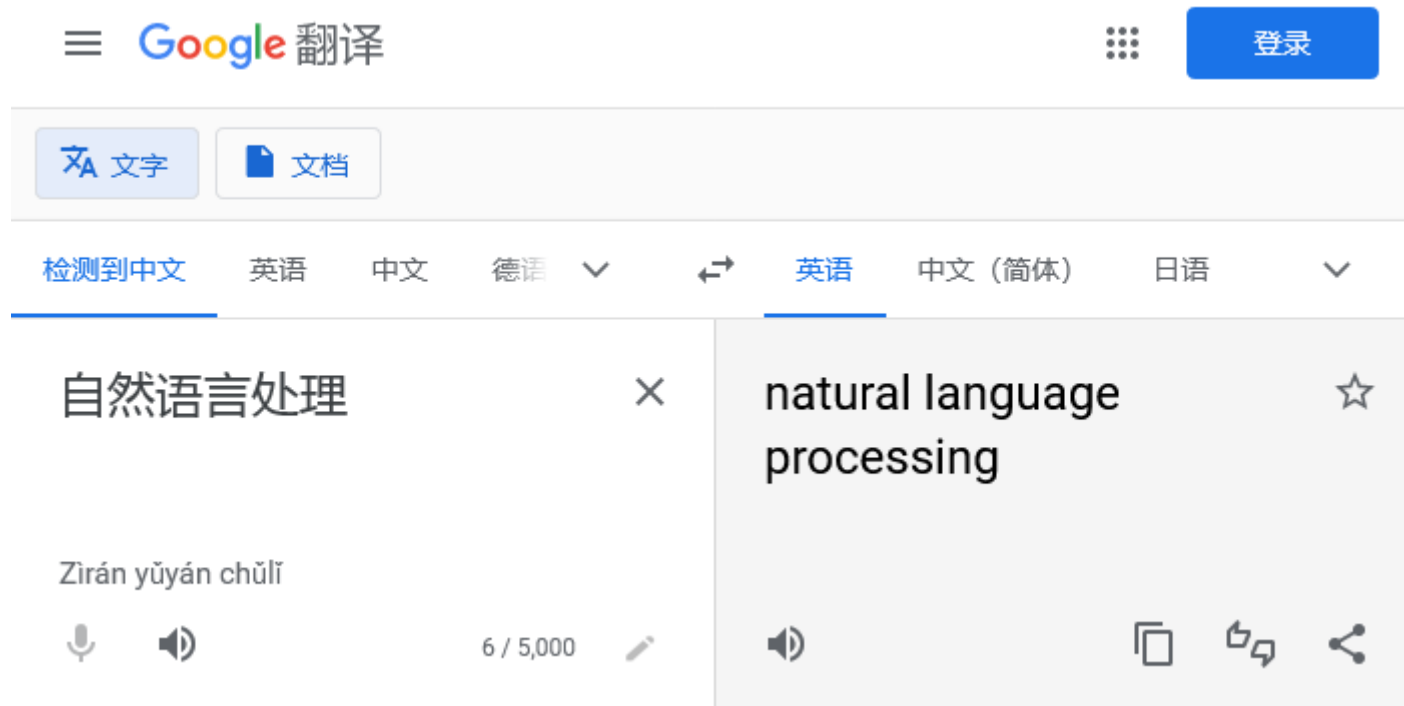


General-Purpose



NLP Applications

► Machine translation



NLP Applications

- ▶ Information extraction
 - ▶ Financial and law documents
 - ▶ E-commerce

← 收件人地址填写

...

📍

粘贴地址信息，自动拆分姓名、电话和地址

📷

收

收件人

📁 地址簿

姓名

电话

- 分机号

城市 / 区域

▼

详细地址（例如：**街**号**）

📍

公司名称（选填）

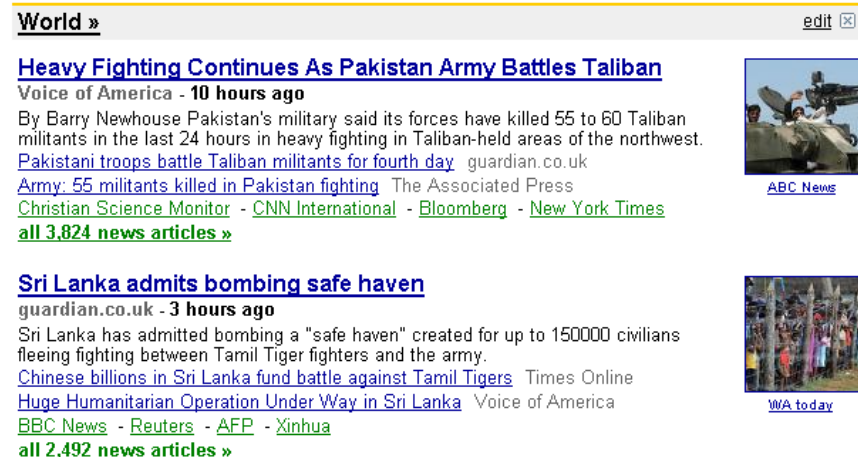
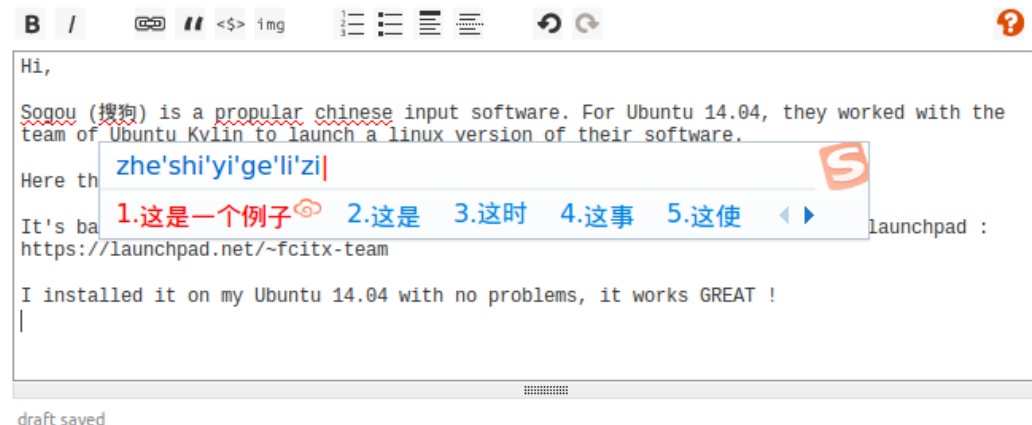
☒ 保存到地址簿

清空当前信息



NLP Applications

- ▶ Chinese IME
- ▶ Grammatical checker
- ▶ News clustering
- ▶ Summarization
- ▶ News generation
 - ▶ Stock market, sports, ...



NLP Applications

- ▶ Essay scoring
 - ▶ Used to score TOEFL and GRE tests!



About the *e-rater*® Scoring Engine

What Is the *e-rater*® Engine?

The *e-rater* engine is an ETS capability that identifies features related to writing proficiency in student essays so they can be used for scoring and feedback. Among other applications, the *e-rater* engine is used within the [Criterion® Online Writing Evaluation Service](#).

Feedback



NLP Applications

- ▶ Art creation (poem, lyric, novel, ...)

藏头诗

九歌

律诗

五言藏头

七言藏头

较悲伤

中性

较喜悦

生成结果

语不随流水

言因有钓矶

处来无限意

理道是吾非



Please write a poem about natural language processing.



A world of words, where machines come alive,
With algorithms that can decipher and thrive,
A language so complex, made simple and clear,
Through the power of NLP, we have nothing to fear.

The voice of a computer, no longer robotic,
But human-like, a new way to be proactive,
With sentiment analysis, emotions at play,
A better way of communication, every day.

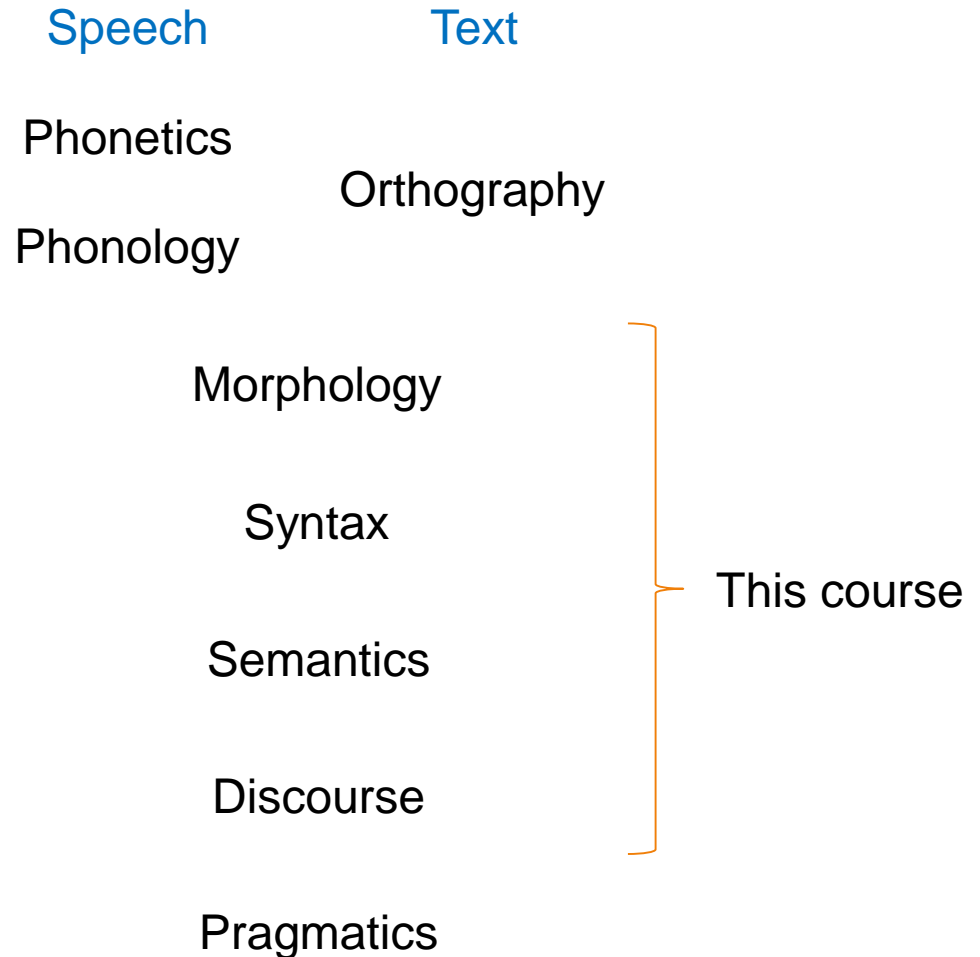
Text summarization, a new way to read,
The essence of meaning, no longer in need,
Of hours of reading, or endless pages to scan,
NLP does it all, in a matter of seconds, it can.

Big data and machine learning, working hand in hand,
To revolutionize the way we understand,
The complexities of language, so much more to come,
Thanks to NLP, we've only just begun.

So here's to NLP, a language so bright,
A new way of communicating, it's a wondrous sight,
A tool of the future, helping us all,
To unlock the power of language, and have a ball!

Why is NLP hard?

- ▶ Human language is complicated!
 - ▶ Levels of linguistic studies:



Why is NLP hard?

- ▶ Language processing requires many levels of knowledge
- ▶ Positive or negative?
 - ▶ The burger tastes bad. *word meaning*
 - ▶ The burger does not taste good. *syntax*
 - ▶ I would not say that the burger is not good. *pragmatics*
 - ▶ “The drink is great!”
“How about the burger?”
“Well...”
 - ▶ The burger tastes like fast food. *world knowledge*



Why is NLP hard?

- ▶ Language processing requires many levels of knowledge

A ship-shipping ship, shipping shipping-ships.



word meaning
morphology
syntax
world knowledge



Why is NLP hard?

- ▶ Ambiguity!

- ▶ Word meaning

- ▶ Iraqi Head Seeks Arms
 - ▶ Stolen Painting Found by Tree

- ▶ Syntactic structure

- ▶ Enraged Cow Injures Farmer with Ax

- ▶ Word meaning + syntax

- ▶ Teacher Strikes Idle Kids



Why is NLP hard?

- ▶ Ambiguity!

- ▶ Semantic structure

- ▶ The detective told his assistant: “Every fifteen seconds a cat in this country gives birth...
 - ▶ ...Our job is to find this cat, and stop her!”

- ▶ Discourse

- ▶ The cat doesn't fit in the box because it is too small.
 - ▶ The cat doesn't fit in the box because it is too large.



Why is NLP hard?

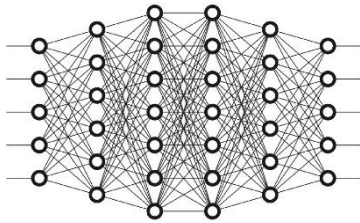
- ▶ Common challenges faced by AI research
 - ▶ High accuracy
 - ▶ Noisy input
 - ▶ Scarce data
 - ▶ Latent variables
 - ▶ Computational efficiency on both space and time
 - ▶ Generalizability
 - ▶ Formal guarantees
 - ▶ Interpretability



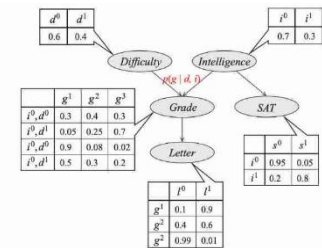
NLP Methodology

Symbolism

$+$ $-$ \times \div
 \neg \vee \perp \approx
 \in \cap \subseteq Σ
 ∂ ∇ \wedge Π



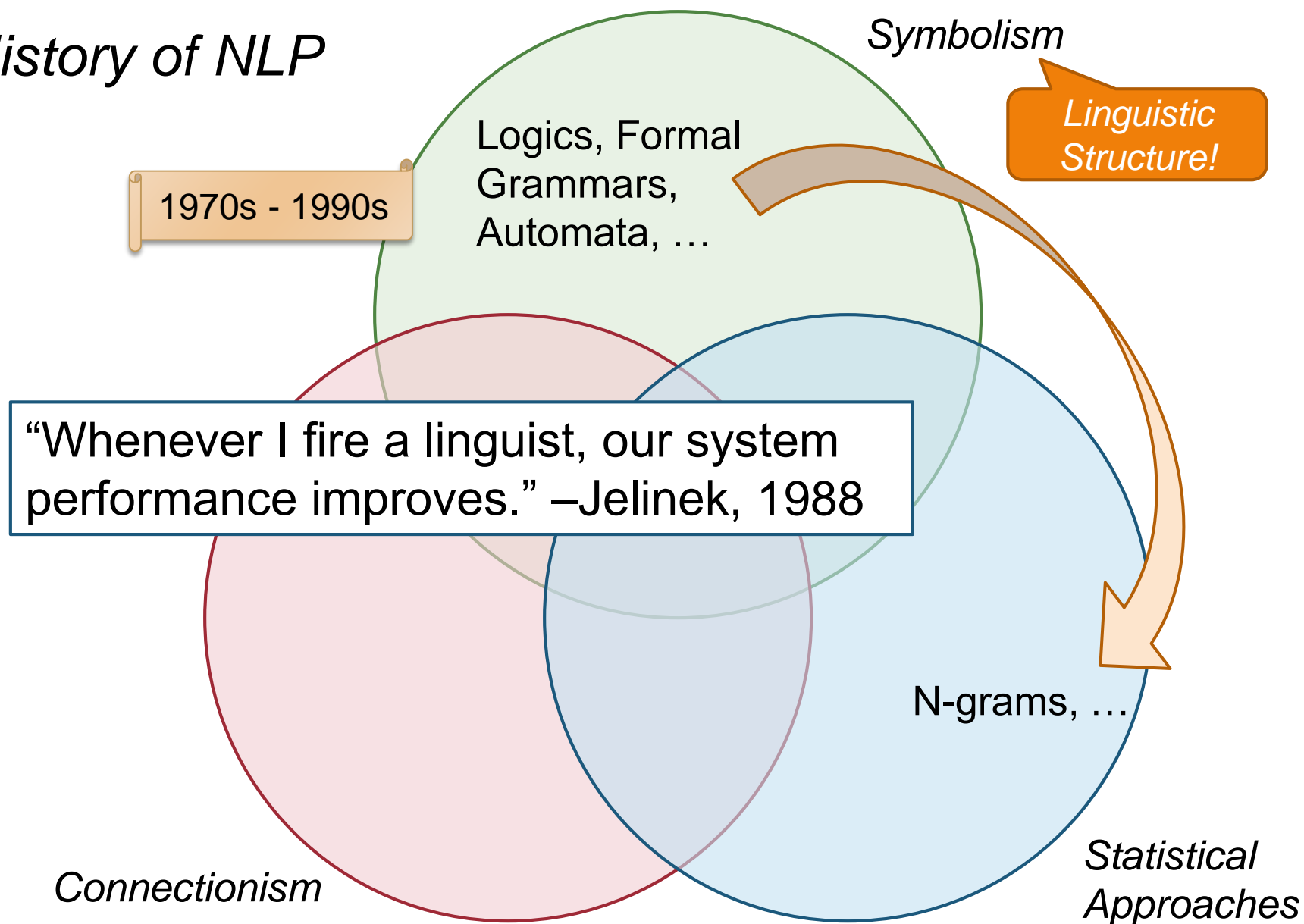
Connectionism



Statistical Approaches

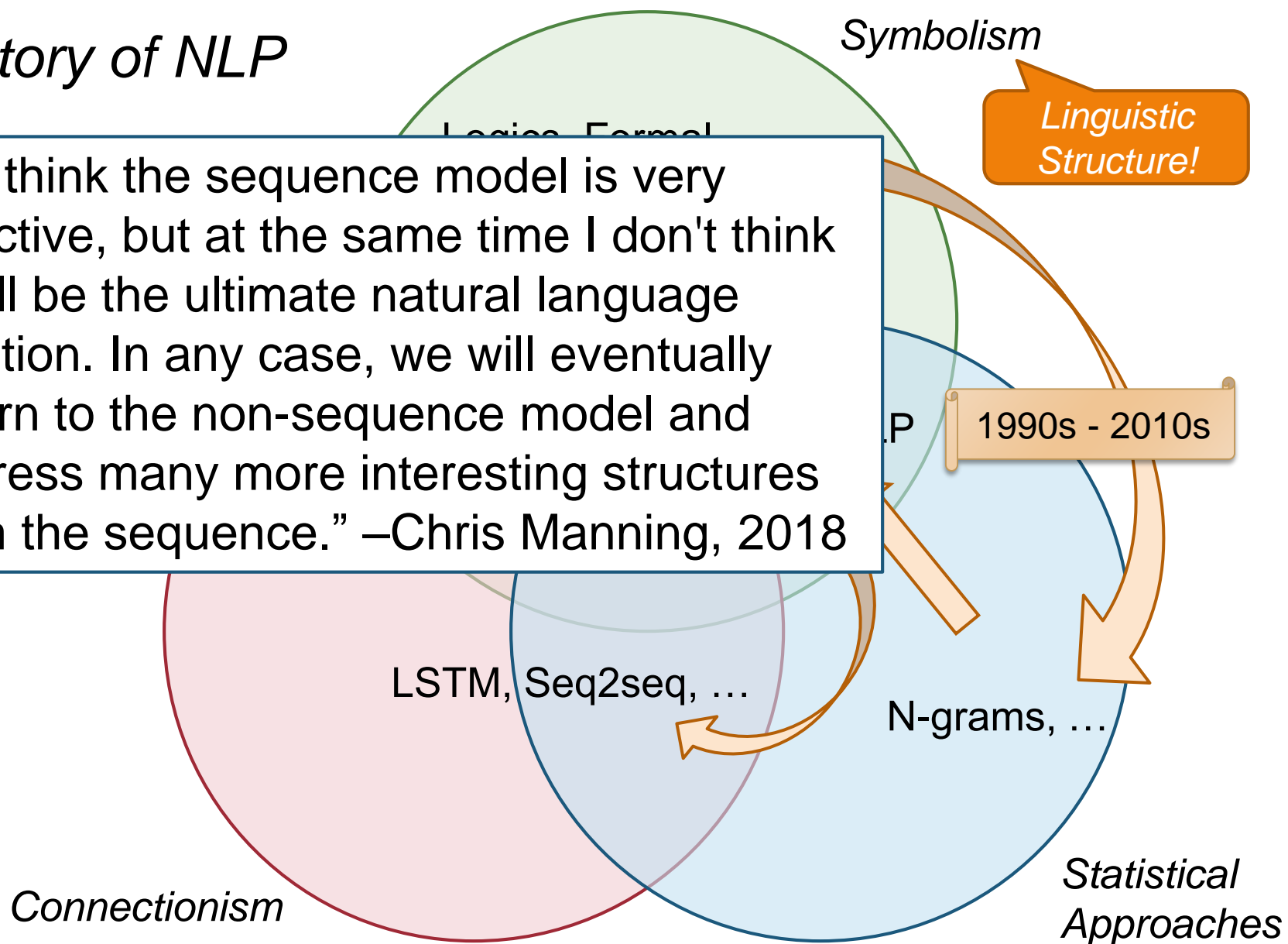


History of NLP

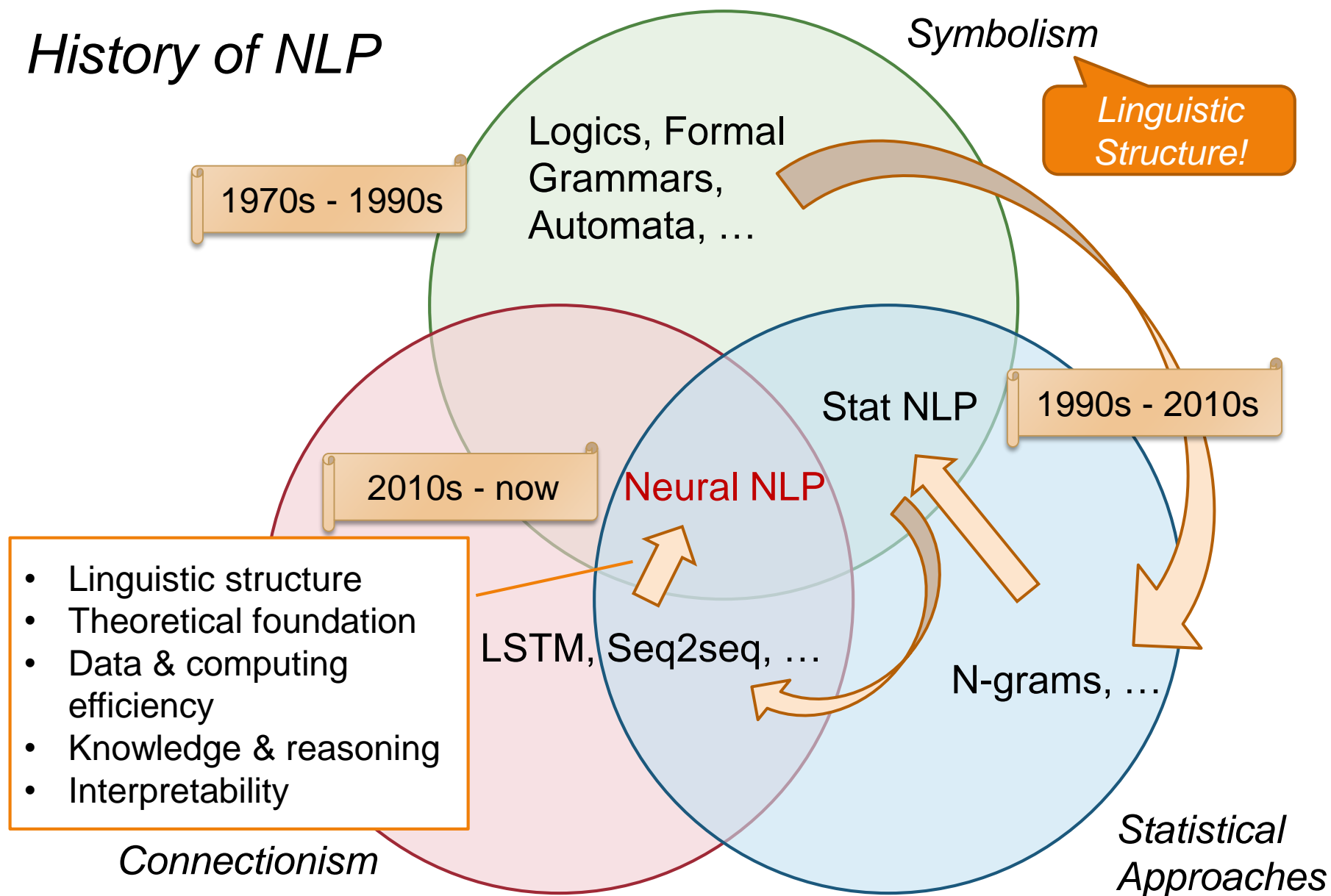


History of NLP

“...I think the sequence model is very effective, but at the same time I don't think it will be the ultimate natural language solution. In any case, we will eventually return to the non-sequence model and express many more interesting structures than the sequence.” –Chris Manning, 2018



History of NLP



Course overview

1. Basics

- ▶ Text normalization
- ▶ Text representation
- ▶ Text classification
- ▶ Text clustering

2. Sequences

- ▶ Language modeling
- ▶ Sequence to sequence
- ▶ Pretrained language models
- ▶ Sequence labeling

3. Structures

- ▶ Constituency parsing
- ▶ Dependency parsing
- ▶ Semantic analysis
- ▶ Discourse analysis

