

# Multimodal LLMs: Speech and Real-Time

Peter Polák, Dominik Macháček  
2/5/2024



Charles University  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

# After the class, you should be able to:

- Know motivation for speech in LLMs
- Know the basic and example speech-to-text methods
- Know real-time methods

## Class outline

- Speech NLP tasks (ASR, translation, emotion recognition, ...)
- Speech in NNs (sound representation, MFCC, raw audio) and in LLMs (Wav2vec, HuBERT, Whisper)
- Simultaneous methods: re-translation vs. incremental
- Streaming policies wait-k and LocalAgreement
- Whisper-Streaming and ELITR demo

# Outline

- Motivation
- Speech and NNs
  - Representation
  - Speech-to-text generation
- Speech and LLMs
- Simultaneous Methods
- Demo

# Speech and LLMs

# Why speech?

Why should LLMs work with speech ... from **user perspective**?

- More natural interaction
- Accessibility
- Some languages do not have a written form
- Complementary with text
- New applications














# Why speech?

Why should LLMs work with speech ... from **technical perspective**?

- Using information beyond text
  - Prosody
    - Intonation, stress, and rhythm
  - Non-verbal language
  - Sentiment
  - Environment
- Avoiding error propagation

# Speech and NNs

# Speech NLP Tasks

- Speech-to-text
  - Automatic Speech Recognition (ASR)  → 
  - Speech Translation  → 
  - Summarization  → 
- Emotion Recognition  → 
- Speaker Verification  +  → 
- Speaker Identification  → 
- And many more ...

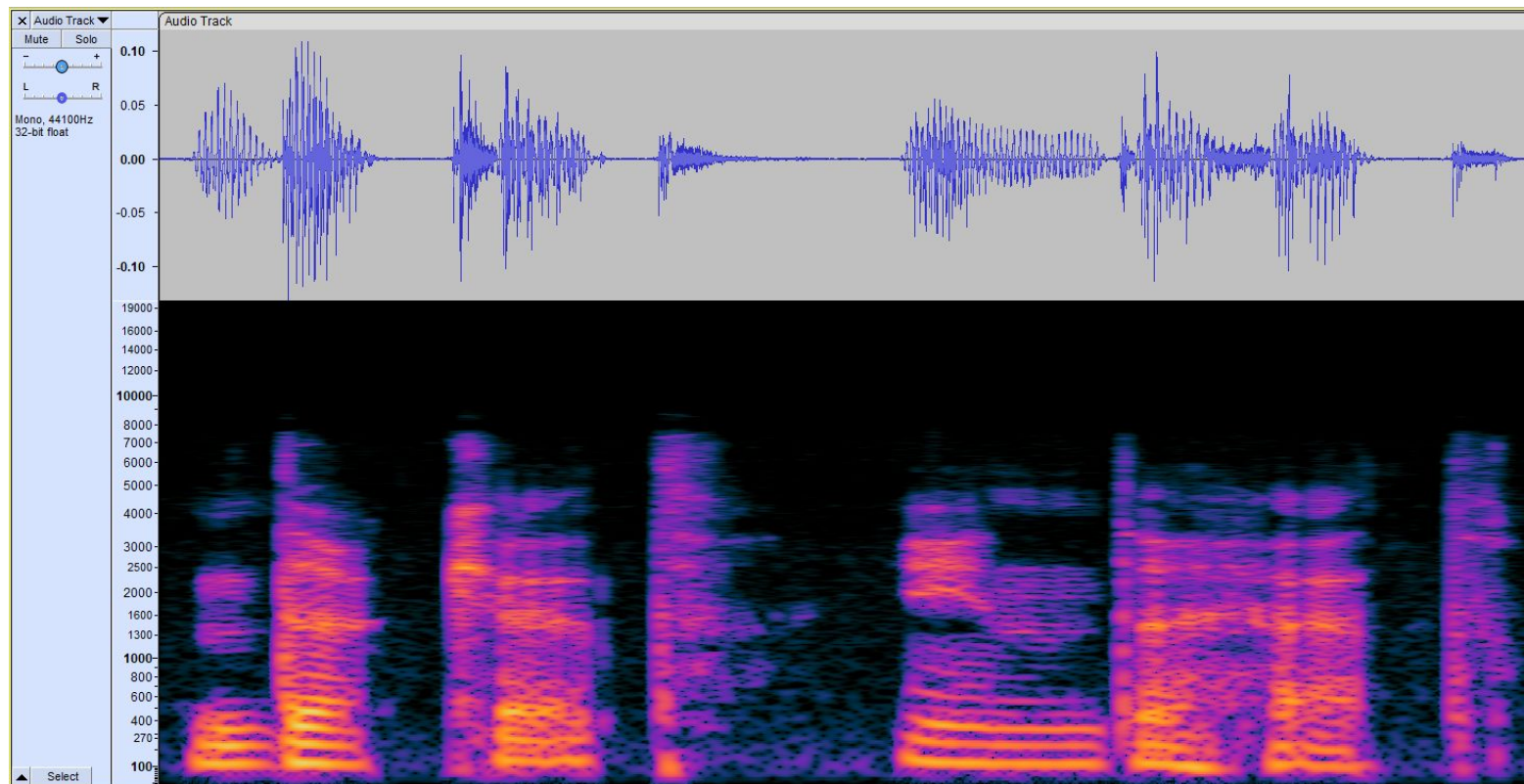


# Speech Representation for NNs

# Speech Representations: Recap from Text

- How neural networks read?
- Tokenization
  - Break the text into smaller units = tokens [characters, words, sub-words]
- Translate tokens to indices in a vocabulary
- Embedding Layer
  - Translate each index into a vector

# Speech Representation



the

cat

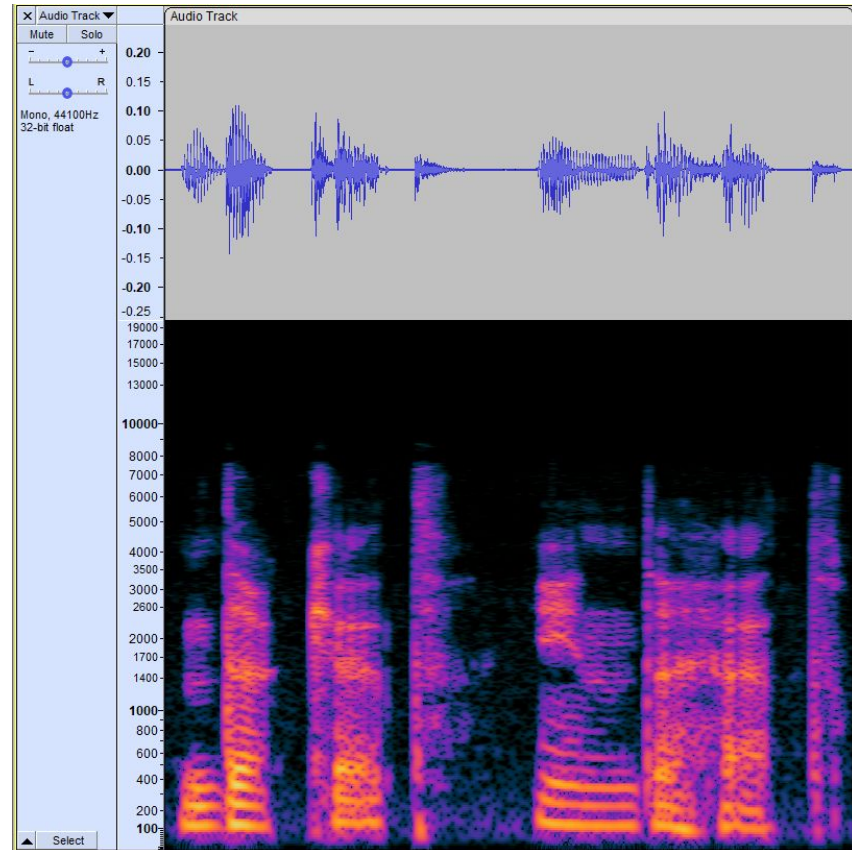
in

the

hat

# Speech Representations: Recap from Text

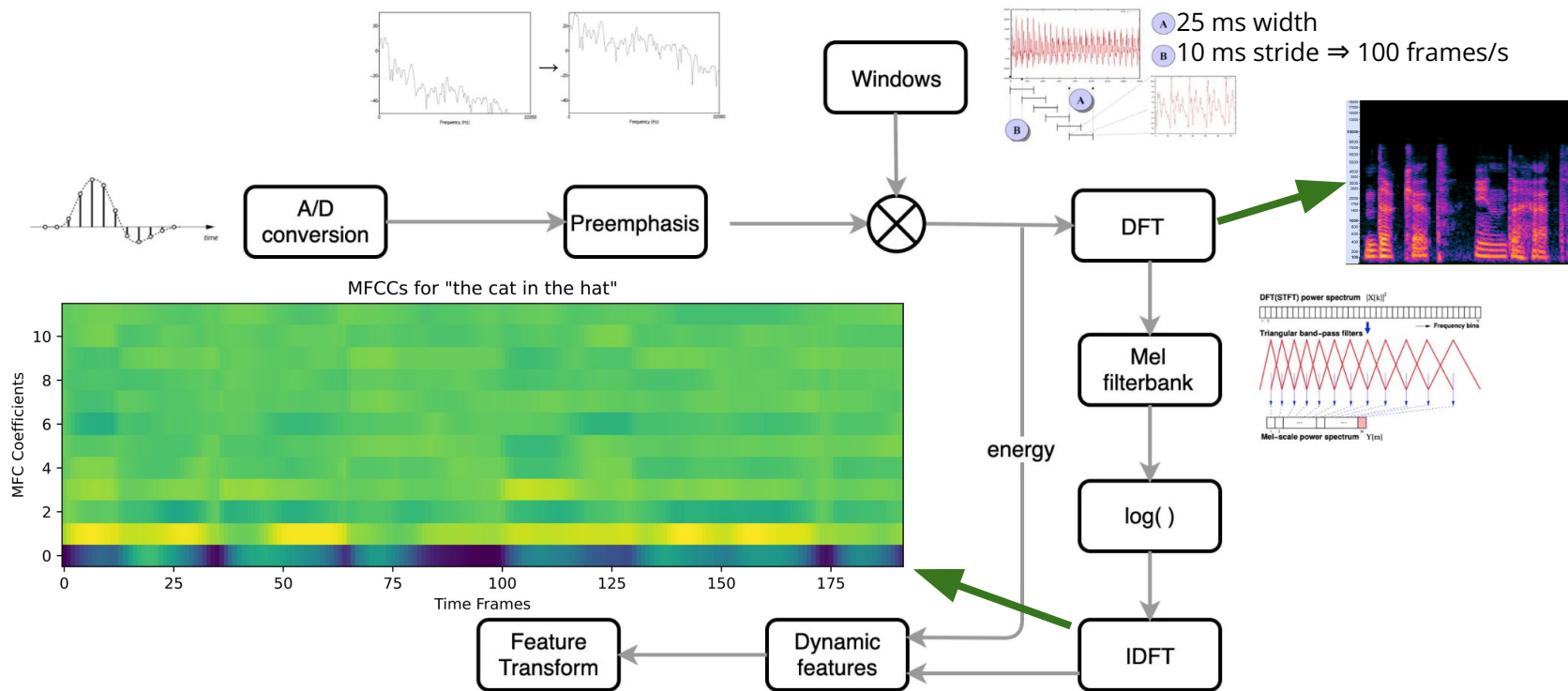
- How do we represent sound in computer?
  - Human speech 100 - 4 kHz, better to 8 kHz
  - 16 kHz wav = 16k floats/s
- How NNs understand speech?
  - Two approaches:
    - MFCCs
    - Raw audio (some tricks needed)



the cat in the hat

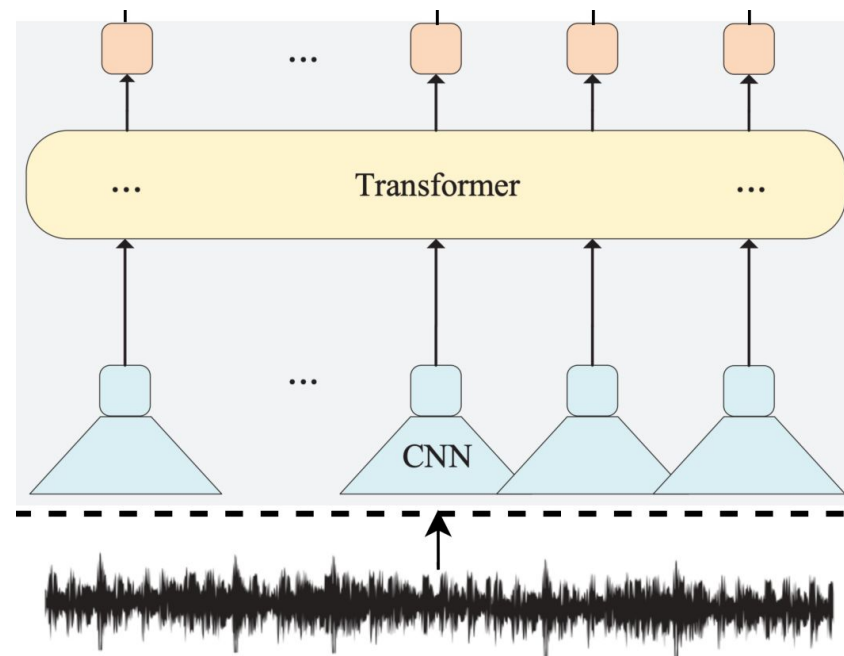
# Speech Representations: MFCCs

- Mel Frequency Cepstrum Coefficients



# Speech Representations: Direct Approach

- Feed directly 16 kHz to NN
  - We want to use Transformers - where is the problem?
    - Complexity of self-attention -  $O(n^2)$
  - Solution
    - Downsample long input with CNNs
    - CNN serves as feature encoder
      - Similar to MFCCs
      - But the representation is **learned from data**
  - Typically a part of pre-trained models
    - Wav2vec 2.0, HuBERT, WavLM



From: SPEECH EMOTION DIARIZATION: WHICH EMOTION APPEARS WHEN?

# Speech Representations: Comparison

- MFCCs

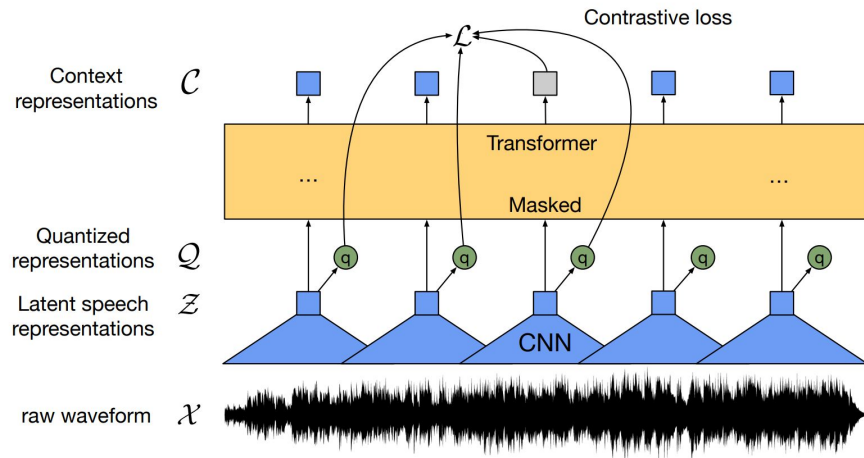
- 👍 Efficient
- 👍 Interpretable
- 👎 Limited features
  - Ideal for ASR, not ER
- 👎 Not robust to noise

- Direct approach

- 👍 More complex features
  - speaker characteristics, emotions, and environmental noise
- 🤔 Depends on training data
  - Can be robust
- 👎 Computational cost
- 👎 Interpretability

# Leveraging Unlabeled Data: Self-Supervised Learning

- We have a lot of unlabeled data
  - Can we use them?
- Pseudo-labeling
  - Use existing ASR
  - 🤔 ASR quality
  - 👎 No ASR for some languages
- Self-supervised learning (SSL)
  - Wav2Vec 2.0, HuBERT, WavLM, ...
- Idea:
  - Model decides on representations based on data
  - Improved context representation using Transformer

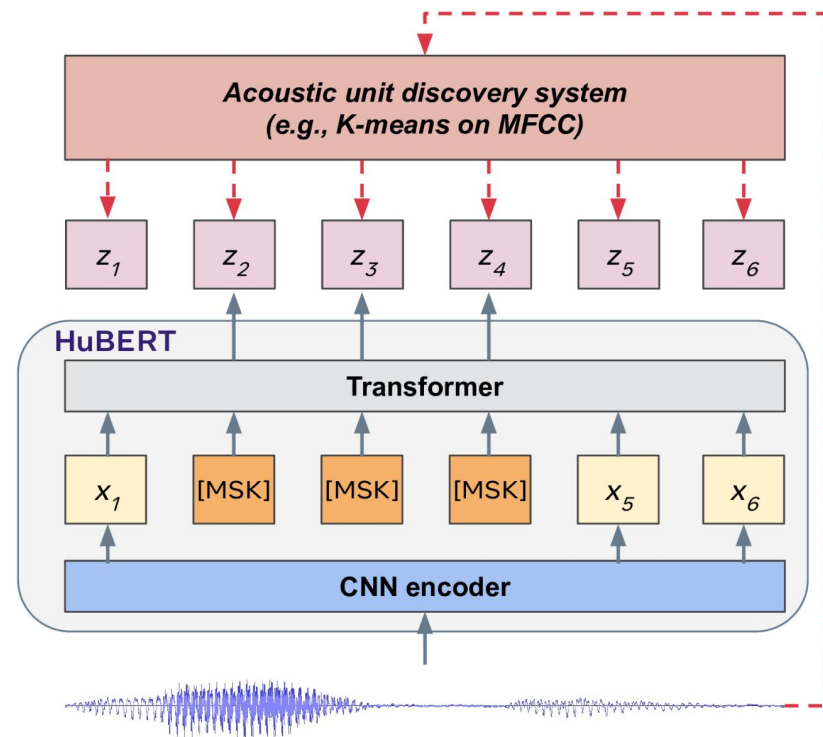


From: wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations



# Leveraging Unlabeled Data: Example

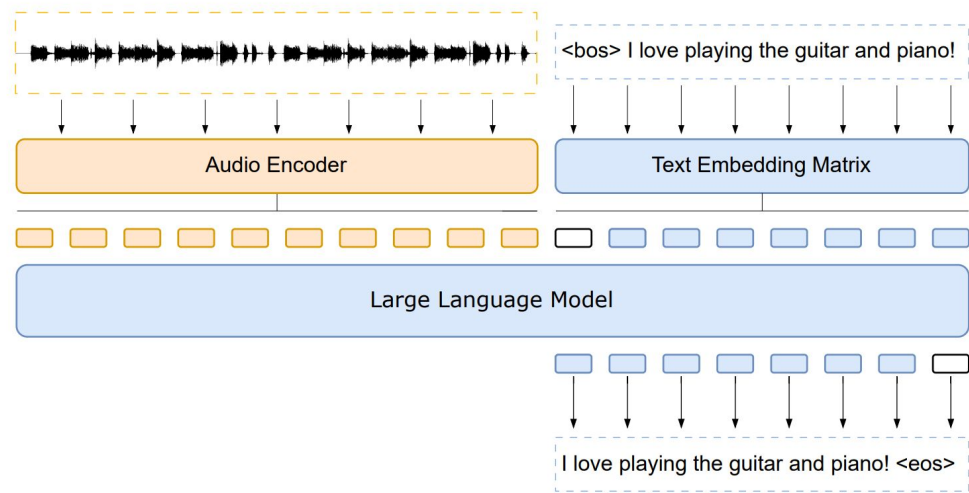
- HuBERT
  - CNN feature encoder
    - 320x downsampling = 50 frames/s
  - Transformer context encoder
- Training
  - Mask random spans
  - Predict acoustic units in masked regions
  - Several hours on up to 256 GPUs
- Acoustic units
  1. Unlabeled audio → MFCCs → k-means (k=100)
  2. Contextual features from Transformer from the first step → k-means (k=500)
- Finetune for downstream task
  - ASR, ER, ...
  - Works even with 10 mins of data



From: wav2vec 2.0: HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units

# Plugging Speech to LLMs

- Hot research topic
- Typically
  - Embed audio with some encoder
    - HuBERT, ...
  - Interleave audio and prompts
  - Feed to LLM



**Figure 2:** Model architecture. The embedding sequence generated from the audio encoder is directly prepended to the text embeddings sequence. This is directly fed into the decoder-only LLM, tasked with predicting the next token. The LLM can be frozen, adapted with parameter efficient approaches such as LoRA or fully finetuned. This work will investigate the former two.

From: Prompting Large Language Models with Speech Recognition Abilities

# Speech-to-Text Example: Whisper

**Multitask training data (680k hours)**

English transcription

- 🗣️ "Ask not what your country can do for ..."
- 📄 Ask not what your country can do for ...

Any-to-English speech translation

- 🗣️ "El rápido zorro marrón salta sobre ..."
- 📄 The quick brown fox jumps over ...

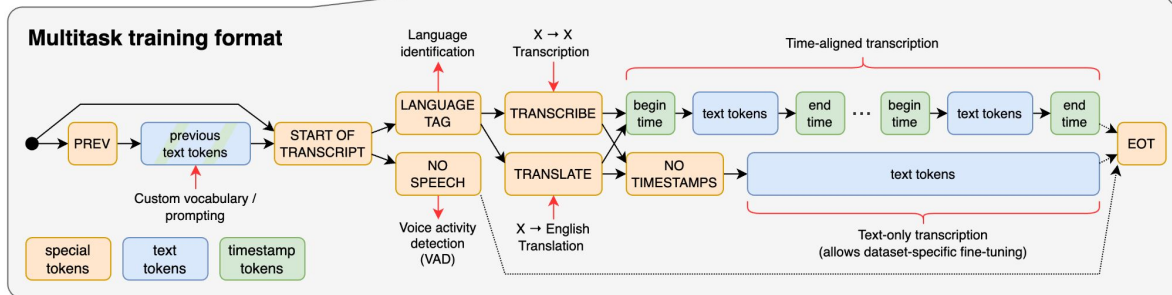
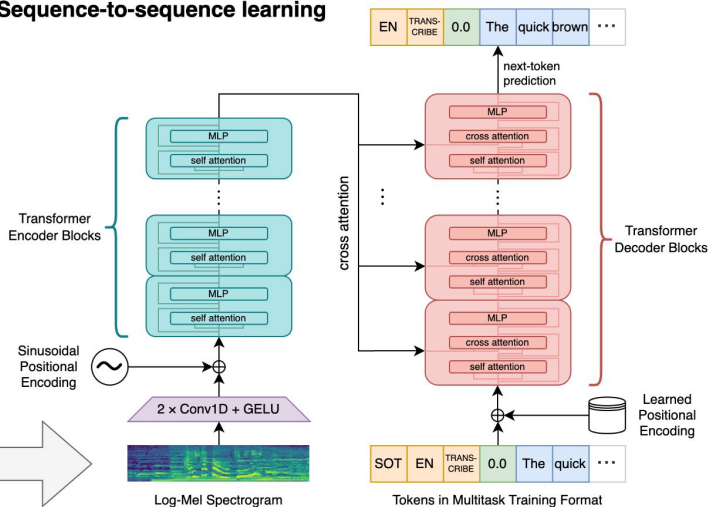
Non-English transcription

- 🗣️ "언덕 위에 올라 내려다보면 너무나 넓고 넓은 ..."
- 📄 언덕 위에 올라 내려다보면 너무나 넓고 넓은 ...

No speech

- 🎧 (background music playing)
- 📄 ∅

## Sequence-to-sequence learning



# LLMs in Real-Time: Simultaneous Speech Translation

Dominik Macháček

2/5/2024, recap from 11/4/2024 + new parts




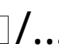







Charles University  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



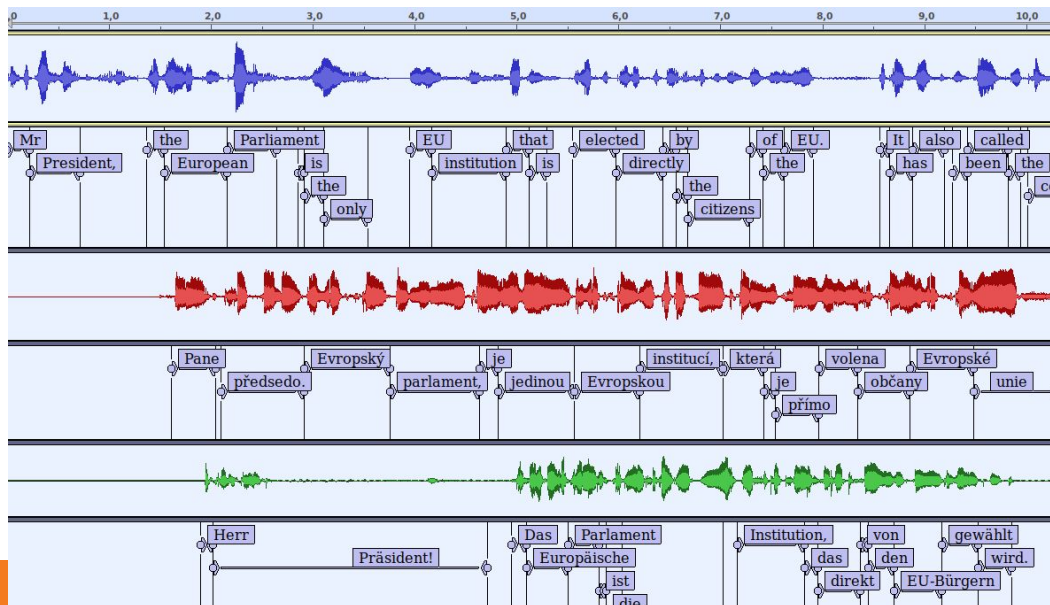
unless otherwise stated

# Problem:

- Seq-to-seq task, e.g.  →  /  /  / ...
- Continuous incremental input    ...
  - = Don't assume <end-of-sequence> mark
  - But expect it in a real-life app. Use e.g. Voice activity detection for 0.5s silence
  
- **Problem:** How to process the task fast?
  - = In **real-time**? 
  - = Simultaneously? 

# Simultaneous speech translation (recap 11/4/2024)

- **Simultaneous** = Live = Real-Time = Low-latency = Incremental
  - Source available continuously, one **chunk** at a time
  - The **chunk** can be:
    - audio segment ... in the direct speech-to-text translation or transcription = ASR
    - or word (text) produced by incremental ASR ... in a cascaded system = ASR + **MT**
  - Provide the target “at the same time” as the source is being produced
    - = simultaneously = with a small additive delay



In the European Parliament:

-> English original source

-> English-to-Czech Sim. Interpreting

-> English-to-German Sim. Intp.

# Simultaneous speech translation (recap 11/4/2024)

- **Simultaneous** = Live = Real-Time = Low-latency = Incremental
  - Source available continuously, one **chunk** at a time
  - The **chunk** can be:
    - audio segment ... in the direct speech-to-text translation or transcription = ASR
    - or word (text) produced by incremental ASR ... in a cascaded system = ASR + **MT**
  - Provide the target "at the same time" as the source is being produced

= simultaneously = with a small additive **delay**



In the European Parliament:

-> English original source

-> English-to-Czech Sim. Interpreting

-> English-to-German Sim. Intp.

# Simultaneous approaches



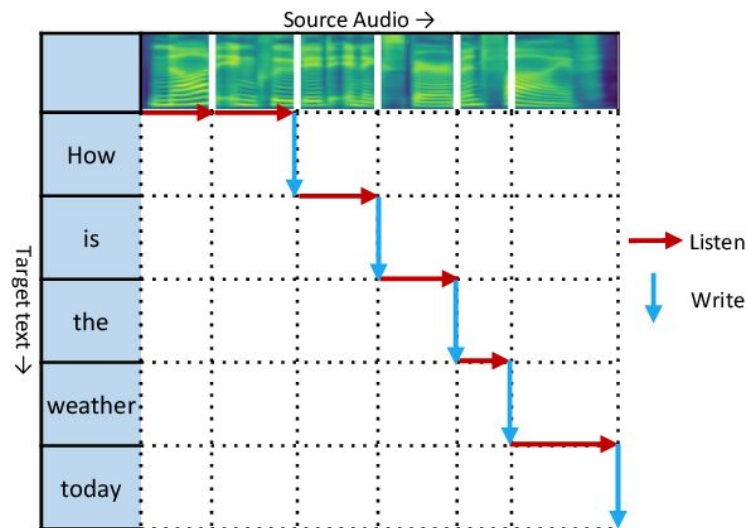
# Re-Translation vs. Incremental

- Re-translate from beginning of sentence each time: **rewrite + append**

Source	Output	Erasure
1: Neue	New	-
2: Arzneimittel	New Medicines	0
3: könnten	New Medicines	0
4: Lungen-	New drugs may be lung	1
5: und	New drugs could be lung and	3
6: Eierstockkrebs	New drugs may be lung and ovarian cancer	4
7: verlangsamen	New drugs may slow lung and ovarian cancer	5
Content Delay	1 4 6 7 7 7 7 7	

Source: [\[Arivazhagan et al., 2020\]](#)

- Alternates between reading from ASR and translating: **no rewrites, only append**



Source: [\[Ren et al., 2020\]](#)

# Re-Translation vs. Incremental

	Re-translation	Incremental
<b>Latency</b>	👍	👎
<b>Quality</b>	👍 as in offline mode	👎 lower than in offline mode
<b>Stability (flicker)</b>	👎 may be unreadable	👍 only appends to the end
<b>Presentation space</b>	Needs large space 👎	Suitable for 2-line subtitles 👍
<b>Versatility</b>	Can serve many users at once 👍 Medium src. lang. proficiency ... needs fast outputs No src. Lang. proficiency ... slower OK, high quality needed	Compromise 👎 latency and quality
<b>Easy to plug-in any offline model?</b>	👍 YES. May be unstable on prefixes but high quality in the end.	Not easy but possible.

# Stability in Re-Translation

# How to make re-translation more stable?

## Baseline (“standard” offline MT)

Source	Output	Erasures
1: Neue	New	-
2: Arzneimittel	New Medicines	0
3: könnten	New Medicines	0
4: Lungen-	New drugs may be lung	1
5: und	New drugs could be lung and	3
6: Eierstockkrebs	New drugs may be lung and ovarian cancer	4
7: verlangsamten	New drugs may slow lung and ovarian cancer	5

**Stability measure:** 13 erasures for 8 generated tokens = **1.625**

# How to make re-translation more stable?

## Baseline (“standard” offline MT)

Source	Output	Erasure
1: Neue	New	-
2: Arzneimittel	New Medicines	0
3: könnten	New Medicines	0
4: Lungen-	New drugs may be lung	1
5: und	New drugs could be lung and	3
6: Eierstockkrebs	New drugs may be lung and ovarian cancer	4
7: verlangsamten	New drugs may slow lung and ovarian cancer	5

Stability measure: 13 erasures for 8 generated tokens = 1.625

## Improvement

Source	Output	Erasure
1: Neue	New	-
2: Arzneimittel	New drugs	0
3: könnten	New drugs may	0
4: Lungen-	New drugs may lung	0
5: und	New drugs may lung and	0
6: Eierstockkrebs	New drugs may lung and ovarian cancer	0
7: verlangsamten	New drugs may slow lung and ovarian cancer	5

4 erasures for 8 generated tokens = 0.5

### We learned: Proportional prefix training

Full	Source	Die Führungskräfte der Republikaner rechtfertigen ihre Politik mit der Notwendigkeit , den Wahlbetrug zu bekämpfen [15 tokens]
	Target	Republican leaders justified their policy by the need to combat electoral fraud [12 tokens]
Prefix	Source	Die Führungskräfte der Republikaner rechtfertigen [5 tokens]
	Target	Republican leaders justified their [4 tokens]

Table 2: An example of proportional prefix training. Each example in the minibatch has a 50% chance to be truncated, in which case, we truncate its source and target to a randomly-selected fraction of their original lengths, 1/3 in this example. No effort is made to ensure that the two halves of the prefix pair are semantically equivalent.

1. Train a standard offline MT
2. Finetune on 1:1 mix of full sent. pairs and src-target prefixes
3. Create the prefixes from the length proportion,
4. do not care about the parallel words in the truncated suffix => anticipation
5. Measure the MT quality and erasures
6. Select a suitable trade-off

**Demo time**

# Demo: Notice the flicker

- From our paper:  
[Presenting Simultaneous Translation in Limited Space, Macháček and Bojar, 2020](#)
- How to present re-translation? **Learn this:**
  - If possible, **use large text output** size, e.g. paragraph view.
    - The end users choose where to look
  - When only few lines are possible:
    - Re-writes can happen in the scrolled away part
    - => **ASR** re-writes usually small ... **OK for few lines**
    - => **MT** re-writes too large  
... too much flicker or large presentation delay, **not usable for few lines**



● On your device: [quest.ms.mff.cuni.cz/elitr/demo/](http://quest.ms.mff.cuni.cz/elitr/demo/)

or QR code →

● Lang. options:

- EN = any lang. translation ... Whisper-Streaming = voice activity controller + auto lang. detection +
- ASR = any lang. ASR + W. large-v3 + long-form streaming with LocalAgreement-2, 1 sec. min. chunk size
- Others\* = En -> 44 langs. MT with re-translation ... UEDIN English-all rainbow model, 44 l. variants, 12-layer enc., 2021-12-10, v4
- \*ZH = now blank, ready to connect another MT

● Wrapped in the **ELITR** pipeline (HE project 2019-2022, CUNI+UEDIN+KIT+others)

<p>powerful experience.</p> <p>80. That was Josef Pаздерка from Czech Radio Plus.</p> <p>81. And now, the President of the Czech Republic, Petr Pavel.</p> <p>82. Please come up here on stage, and present your opening speech to start the first session of this conference.</p> <p>83. Mr. President.</p> <p>84. Good morning, ladies and gentlemen, guests here and listeners and viewers on the other platforms.</p> <p>85. When I was asked by the Czech radio to take over the auspices of this event, I did not hesitate for a second, because the topics that we are discussing here today are very important to me.</p> <p>86. This is the 100th anniversary since the start of the regular broadcast of the Czech radio, which also tells us about the importance of freedom of speech, of talking without censorship, without limitations, the freedom to accept information, to seek information, to spread information, the freedom that in many parts of the world is restricted very strongly, and a freedom... people keep giving their lives for.</p> <p>87. And specific examples are not far away.</p> <p>88. We have among us the daughter of Boris Nemtsov, the murdered Russian opposition politician, Zhanna Nemtsova.</p> <p>89. On Vinohradská street, quite close to the headquarters of the Czech Radio, there is Radio Free Europe, and three of its journalists are now in prison,</p>	<p>terpánky, petra pavla, aby přišel sem k nám a přednesl svůj úvodní projev a vlastně tak otevřel ten první blok celé konference.</p> <p>60. Blok nazvaný Ukrajina jako společná odpovědnost.</p> <p>61. Prosim, pane prezidente.</p> <p>62. Dobry den, damy a panove, vazeeni hoste zde v sale, posluchaci, ale take divaci na ostatnich platformach.</p> <p>63. Kdyz me vedeni Ceskeho rozhlasu pozadalo o zastitu nad dnešní konferencí, nemusel jsem dlouho váhat, protože tématá, kterými se tady zabýváme, jsou pro mě velice důležitá.</p> <p>64. Připomínáme si 100. výročí odzahájení pravidelného rozhlasového vysílání a to je zároveň i připomínkou významu svobody slova.</p> <p>65. Svobodu vyjadřovat se bez cenzury a bez omezení.</p> <p>66. Svobodu přijímat informace a myšlenky, vyhledávat je a šířit.</p> <p>67. Svobodu, která je v různých koutech světa stále výrazně omezována a za její šproszování lidé i dnes platí tu nejvyšší cenu.</p> <p>68. Pro konkrétní příklady nemusíme vůbec chodit daleko.</p> <p>69. Mezi námi je dnes dcera zavražděného ruského opozičního politika Borise Němcova, žena Němcovová.</p> <p>70. Na ulici Vinohradská, jen kousek od sídla Českého rozhlasu, sídlí i Radio Sobotná Evropa.</p> <p>71. Jelikož tři novináři jsou dnes vězněni. – Jiřard Lošak a Andrej Kuzněčik v Bělorusku a Vladislav Jespenko na ruském okupovaném Krymu.</p> <p>72. V únoru tohoto roku jsme si připomněli</p>	 <p>83. is-sur President.</p> <p>84. Filghodu tajieb, nisa u mara, mistiedjina hawn u dawk li jisinghu u Hspetturi fuq il-</p> <p>85. Meta ntablani mir-radju Ċeka biex tiehu l-awditi ta' dan l-aveniment, ma stajfx ghal sekonda, minhabba li s-sugġetti li qed niddiskutu hawn illum huma importanti hafna ghajlia.</p> <p>86. Dan huwa li-100 anniversarju mill-bidu tat-trażmissjoni regolari tar-radju Ċek, li jghidna wkoll dwar l-importanza tal-libertà tal-kunsiderazzjoni, ta' tkellem minghajr cenzura, minghajr limitazzjonijiet. Il-libertà li jaçċettaw informazzjoni, li jiftixu informazzjoni, li jinfirxu informazzjoni, il-libertà li f'hafna partijiet tad-dinja hija ristretta hafna b'saħħitha, u liberta... in-n</p> <p>87. U ezempji specifiċi mhumiex boghod.</p> <p>88. Ahna għandna fostna t-tifla ta' Boris Nemtsov, il-politika ta' l-oppożizzjoni Russa maqulta, Zhanna Nemts</p> <p>89. Fuq id-nejn ta' Vinohradská, appo hafna mill-kwartieri generali tar-Radju Ċeka, hemm ir-Radju Hlelsa ta' l-Ewropa, u</p>	<p>вступіне слово, і відкрив перший блок конференції.</p> <p>70. Блок під назвою «Україна як спільна відповідальність».</p> <p>71. Прошу пана президента.</p> <p>72. Доброго дня, дами та панове, доброго дня гості в залі, слухачі, а також глядачі на інших платформах.</p> <p>73. Коли керівництво Чеського радіо попросило мене взяти патронат на цієї конференції, я не вагався.</p> <p>74. Тому що на теми, про які ми сьогодні будемо говорити, це теми дуже важливі для мене.</p> <p>75. Сьогодні ми пригадуємо соту річницю від початку трансляції Чеського радіо.</p> <p>76. І це також нагадування про важливість свободи слова.</p> <p>77. Свободу висловлювати свою думку без обмежень.</p> <p>78. Свободу приймати інформацію та думки, шукати їх та поширювати.</p> <p>79. Свободу, яка у різних частинах світу досі піддається переслідуванням.</p> <p>80. І за неї люди і сьогодні платять високу ціну.</p> <p>81. За такими прикладами нам не треба ходити далеко.</p> <p>82. Сьогодні серед нас є донька вбитого російського політика Жанна Німцова.</p> <p>83. У вулиці Віноградська, зовсім недалеко від місця, де знаходиться Чеське радіо, знаходиться і радіо «Свобода».</p> <p>84. Три журналіста, якого зараз знаходиться за ґратами.</p> <p>85. У лютому цього року</p>	<p>وموضوع هذا المؤتمر واضح تماما. للتأكيد على نوعية المعلومات التي تأتي إلى الجمهور التشيكي. من المهم ما نوع المعلومات التي يستهلكها.</p> <p>75. وأنيأ لم تزعج، بعد هذا أكثر من سنة الصراع، لم ترغب في أن ينظر إلى هذا الألعاب فيديو.</p> <p>76. بعض الحركات على الخريطة وما زال هناك مصير فرادي الناس، معانا فرادي ما يحدث.</p> <p>78. طوابع هذه الألواع، طوابع اليوم، يجب أن تكون قادراً على رؤية على الأقل نظرة على ذلك.</p> <p>79. وأمل أن تكون تحريه منيرة للأهتمام وقوية. كان ذلك (أخيراً باردركا) من الراديو التشيكي (ر).</p> <p>81. والأن، رئيس الجمهورية التشيكية، بنتر بالي.</p> <p>82. رجاء تعال إلى هنا على المسرح، وعرض حوارك الانتخابي لبدء الدورة الأولى لهذا المؤتمر.</p> <p>83. سيدي الرئيس</p> <p>84. صباح الخير يا سيداتي وسادة صوبف هنا والمستمعين والمساهدين على المنبر الأخرى وعندما طلبت من الأذاعة التشيكية أن تتولى رعاية هذا الحدث، لم أتردد لمدة ثانية، لأن المواضيع التي نناقشها اليوم هامة جداً بالنسبة لي.</p> <p>86. هذا هو الذكرى السنوية المئوية منذ بداية البث المنتظم للأذاعة التشيكية، التي تُحرباً أيضاً بأهميه حرية التعبير، والوحد دون رقابة، دون قيود، حرية قول المعلومات، والوحد على المعلومات، ونشر المعلومات، والحرية التي تُقيد في العبد من أنحاء العالم بقوة جداً، والحرية... الناس يستمرن في إعطاء حياتهم من أجلها.</p> <p>87. وهناك أمثلة محددة ليست بعيدة عن ذلك.</p> <p>88. لدينا ابنة بوبريس نامتسوف، سياسة المعارضة الروسية المغتلة، زاننا نامتسوف.</p> <p>89. في شارع فينوهرادسكا، قريب جداً من مقر الإذاعة التشيكية، هناك إذاعة أوروبا الحرة، واثانة من صحفيا في السجن الآن.</p>
--	---	---	--	---



**Incremental**

# Problem: When to wait or translate

Chinese source:

jǐngfāng      xiàzhōu  
警方      下周  
police      next week

jiāng      duì      bù fèn      shè àn      rén yuán      tí qǐ gōng sù  
将      对      部分      涉案      人员      提起公诉  
will      for      part      involved      people      accuse

Simultaneous intp.:  
Next week, police

will

accuse some of the people involved in the case.

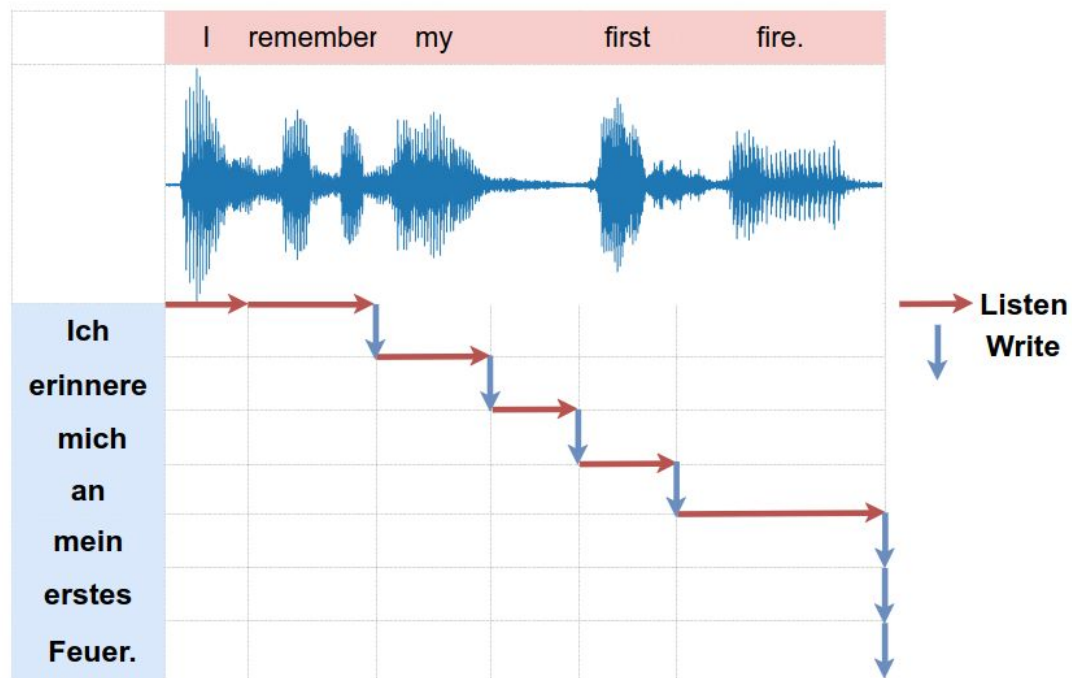
Source: <https://aclanthology.org/2020.emnlp-main.178.pdf>

# Simultaneous policies

Sim. **policy** for an incremental input:

Given a seq-to-seq model, make continuous predictions.

Objectives: latency and quality. Control them by one parameter.



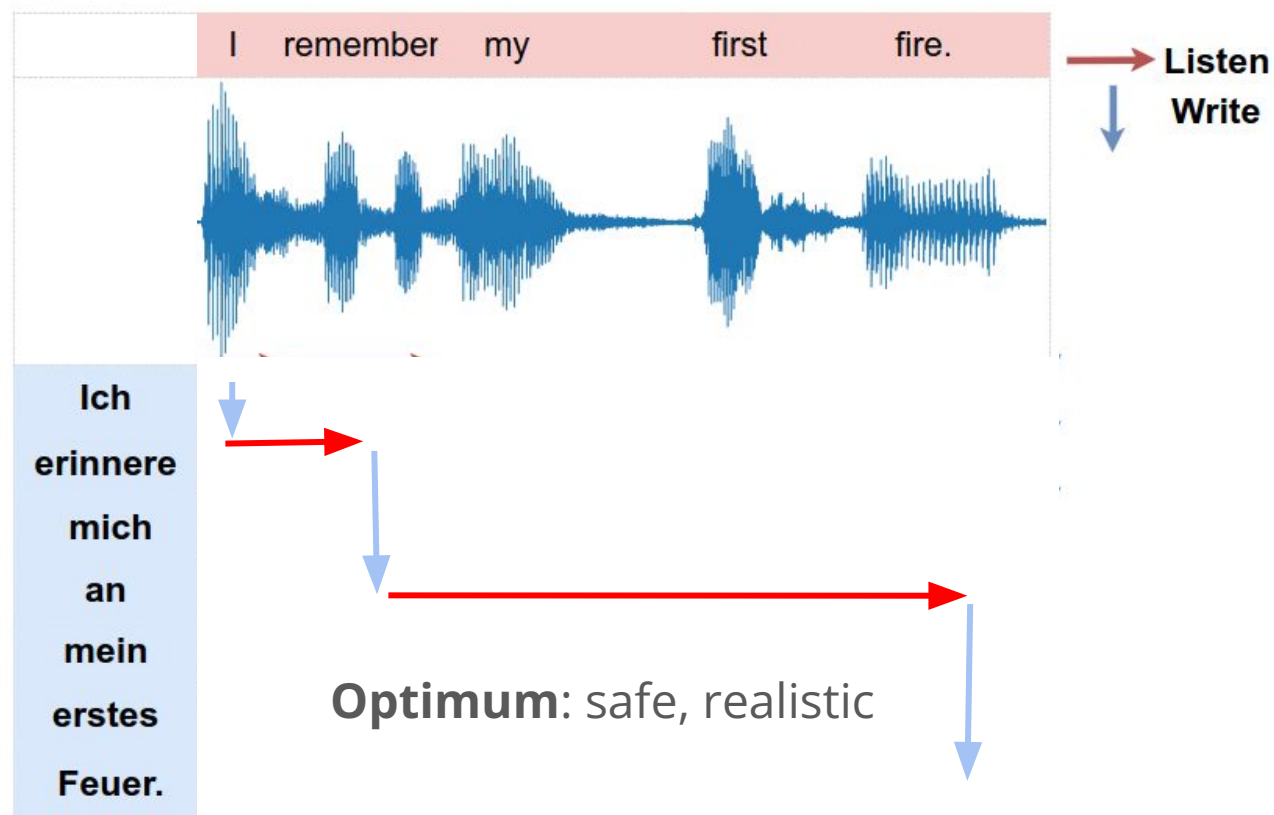
# Simultaneous policies



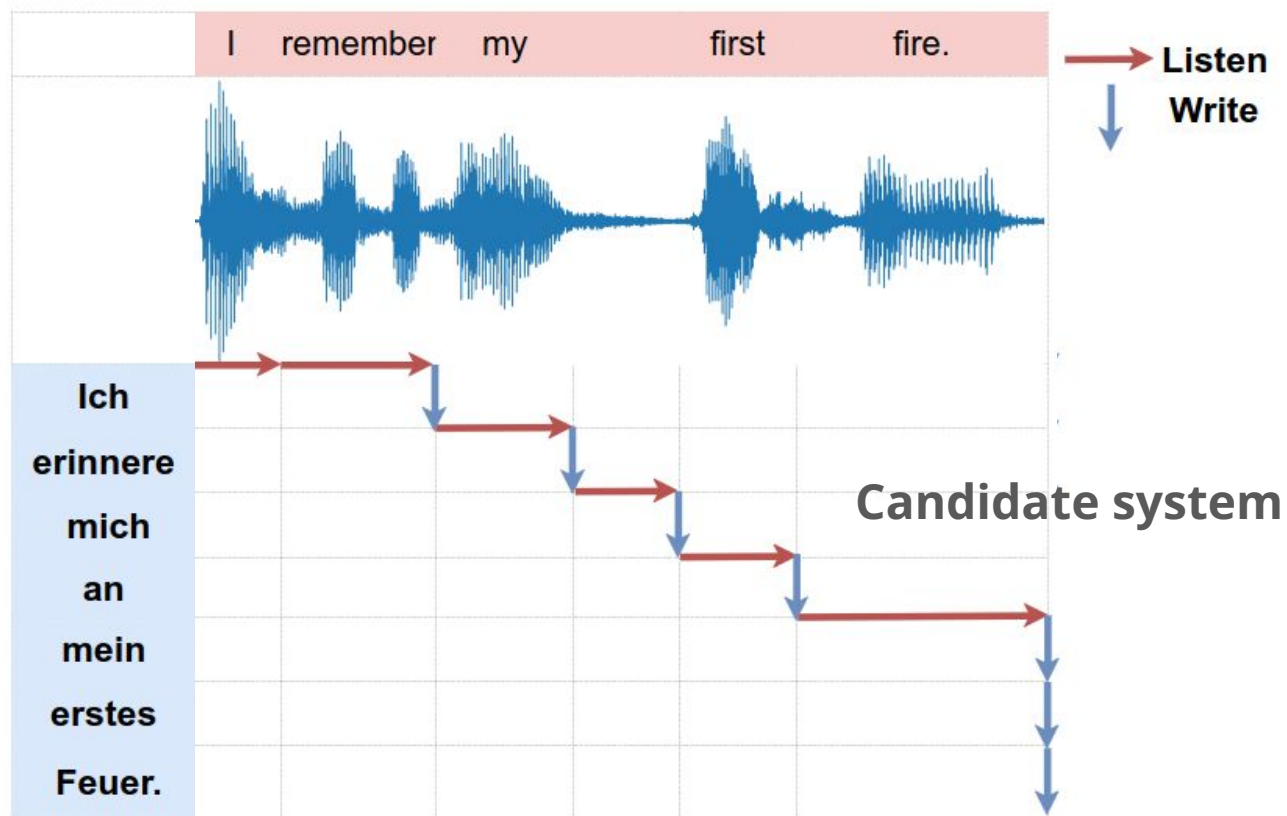
# Simultaneous policies



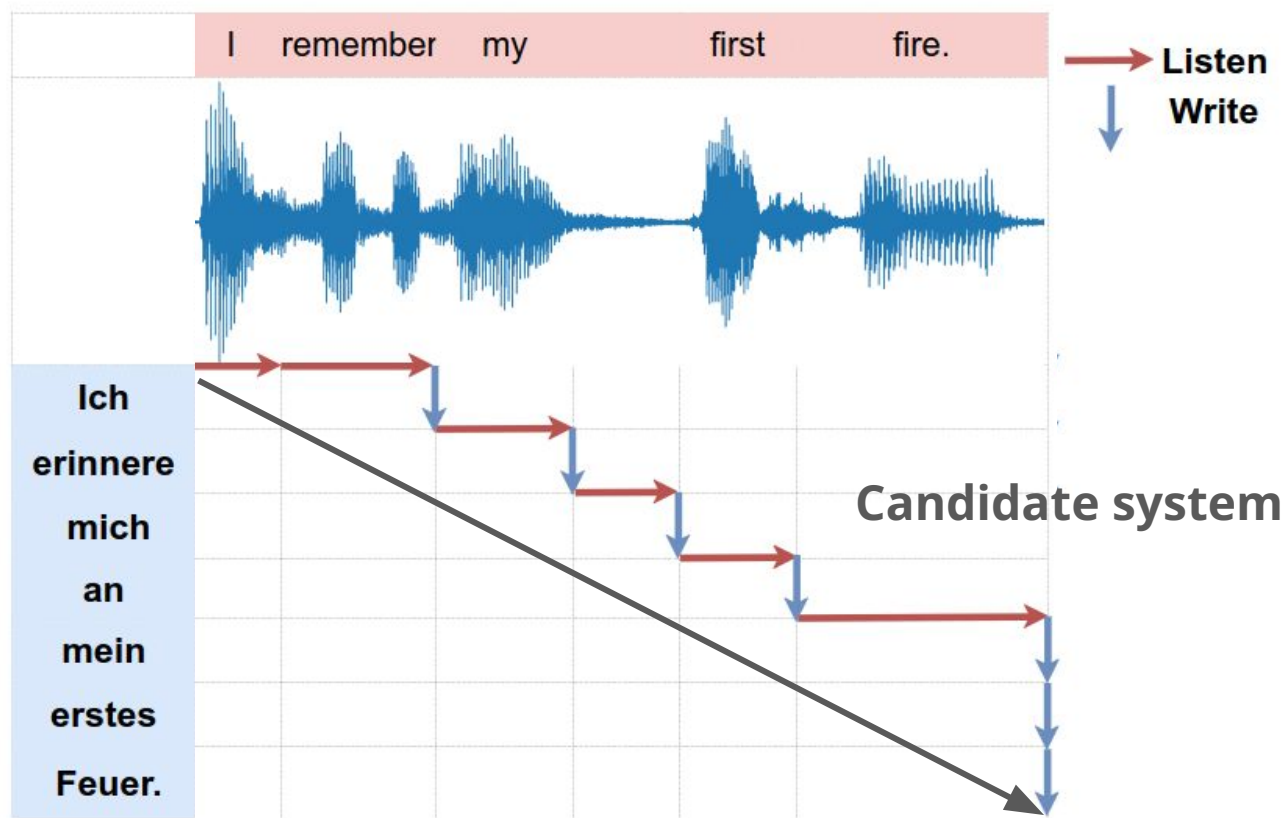
# Simultaneous policies



# Simultaneous policies



# Simultaneous policies

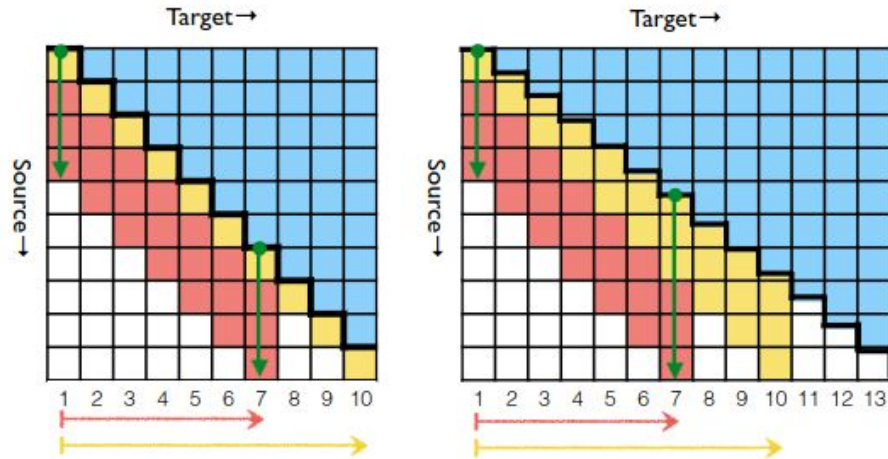


**Diagonal:** Latency reference



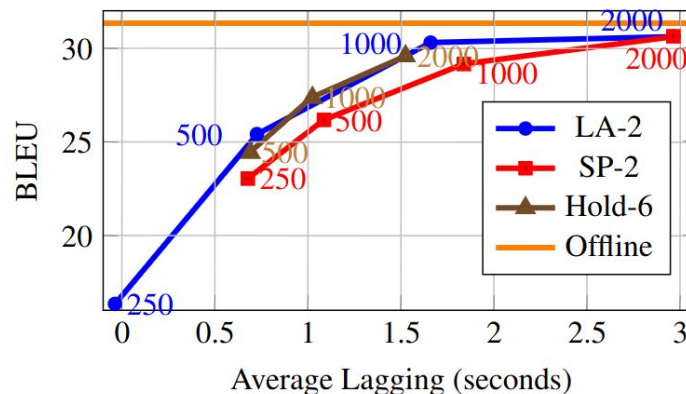
# Simultaneous Policies

- **Wait-k:** Be always **k** words behind.  
(+ catch-up factor because src-tgt length ratio.)
  - Simple baseline
  - Small **k** might be OK for some lang. pairs
  - Not self-adaptable by lang. or content



# Simultaneous Policies

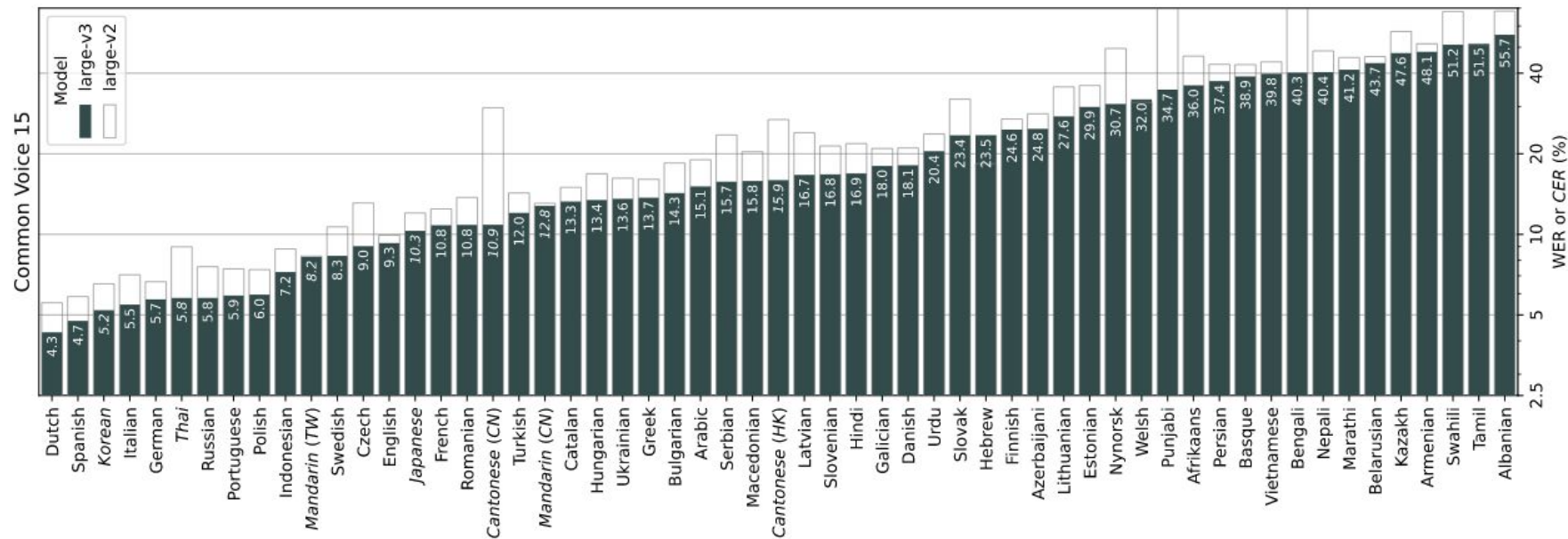
- LocalAgreement-N (Polák et al., 2022):
  - With every new src chunk, decode internally.
  - Output **Agreement** = common prefix of the last **N** chunks.
  - Use it for forced decoding the next chunks.
  - => **Self-adaptable** by lang. and source complexity
  - Parameters: N, the chunk size (e.g. x-times 330ms = appx. x words)
- LocalAgreement-2 (Liu et al., 2020):
  - Best in IWSLT 2022 (Polák et al., 2022)
  - Min. latency 2x chunk size, max unlimited
  - In Whisper-Streaming



- Practical tool, demonstrates SoTA
- Includes:
  - Voice activity controller = end after 0.5s silence
  - + automatic language detection (99 langs)
  - + Whisper, with fast implementation
  - + LocalAgreement-2
  - + heuristics for continuous speech
  - Wrapped in ELITR pipeline

# Whisper Languages

- Let's chat in 99 langs.



# List of Whisper's languages

Afrikaans Albanian Amharic Arabic Armenian Assamese Azerbaijani Bashkir Basque  
Belarusian Bengali Bosnian Breton Bulgarian Burmese Cantonese Castilian Catalan  
Chinese Croatian Czech Danish Dutch English Estonian Faroese Finnish Flemish  
French Galician Georgian German Greek Gujarati Haitian Haitian Creole Hausa  
Hawaiian Hebrew Hindi Hungarian Icelandic Indonesian Italian Japanese Javanese  
Kannada Kazakh Khmer Korean Lao Latin Latvian Letzeburgesch Lingala Lithuanian  
Luxembourgish Macedonian Malagasy Malay Malayalam Maltese Maori Marathi  
Moldavian Moldovan Mongolian Myanmar Nepali Norwegian Nynorsk Occitan  
Panjabi Pashto Persian Polish Portuguese Punjabi Pushto Romanian Russian  
Sanskrit Serbian Shona Sindhi Sinhala Sinhalese Slovak Slovenian Somali Spanish  
Sundanese Swahili Swedish Tagalog Tajik Tamil Tatar **Telugu** Thai Tibetan Turkish  
Turkmen Ukrainian Urdu Uzbek Valencian Vietnamese Welsh Yiddish Yoruba

## Summary

- You learned how to represent speech
- How to leverage unlabeled speech data
- How to integrate speech into LLMs
- Simultaneous methods and policies