# How Has Deep Learning Revolutionized Human Language Technology?

*Tao Ma, Ph.D.*

Vice President of the Language Technology Group
ASAPP Inc.
12/6/2020

**ASAPP**

# Human Language Communication

- **Communication is transferring information, through encoding/decoding/channels.**

- **Human language communication deals with information on a conceptual level, integrated with knowledge of speech production, perception and linguistics.**



ASAPP

# Human Language Technology

How computer can analyze, produce, modify or respond to human texts and speech.

### Speech Recognition

Automatic Speech Recognition (ASR) converts spoken words into text. It detects spoken sounds and recognize them as words.
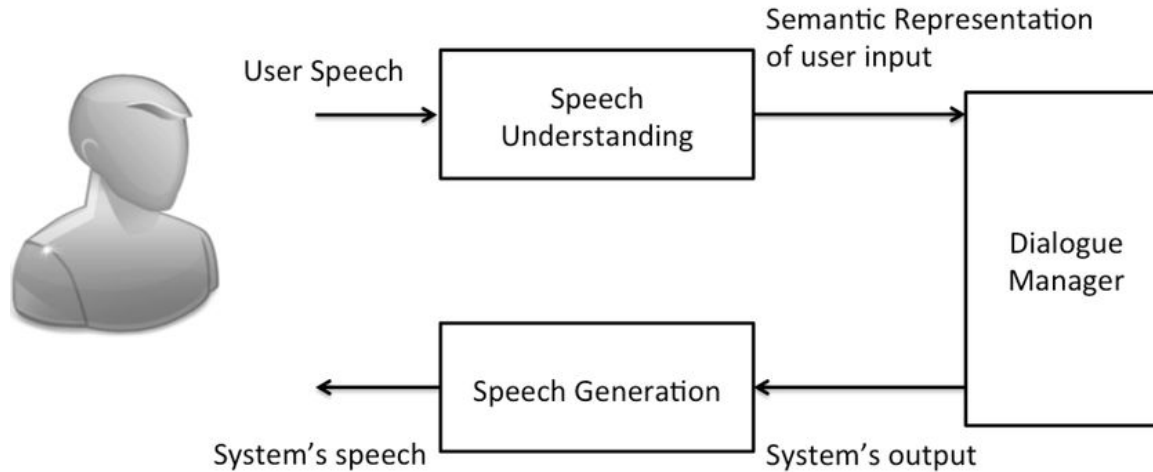
### Speech Synthesis

Text To Speech (TTS) converts text to voice. It's a technology that makes computer talk.

### Natural Language Processing

Natural Language Processing (NLP) helps computers understand, interpret and manipulate human language.
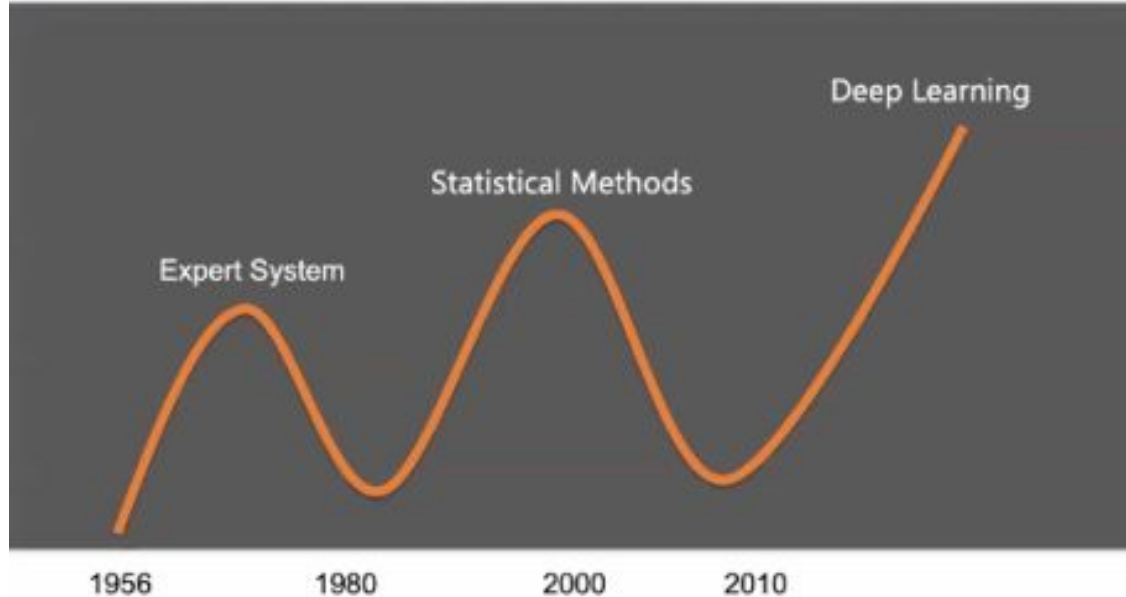
**ASAPP**

# Behind Siri/Alexa/Google Assistant/Cortana

How AI assistants work? The architecture of a Spoken Dialogue System.
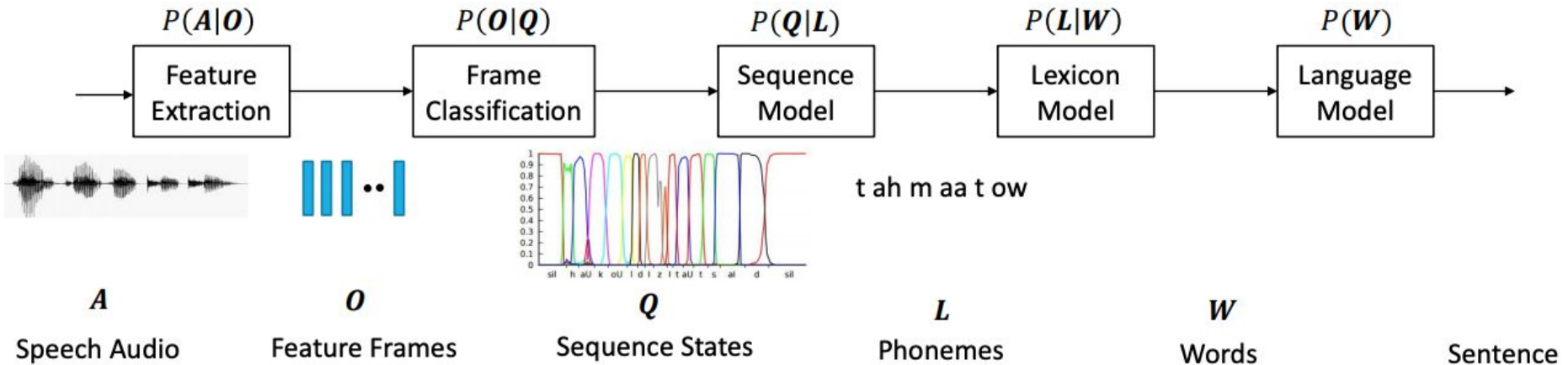
# Deep Learning: 3rd wave of Artificial Intelligence

In 2011, *AlexNet* achieved big wins in Computer Vision, then deep learning is quickly being adopted in Speech and NLP.
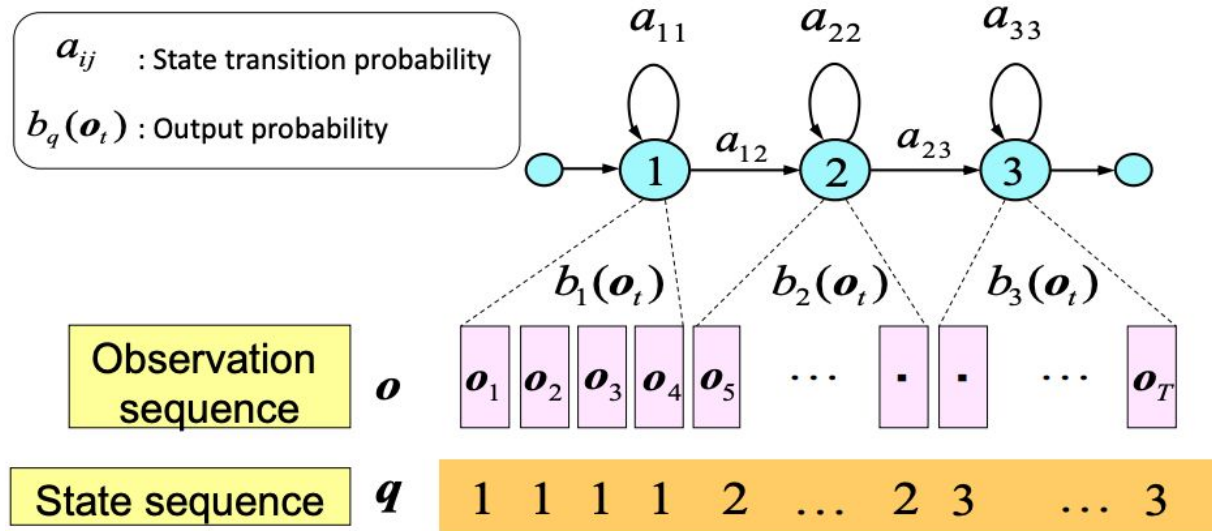
# Speech Recognition

$$\widehat{W} = \operatorname*{argmax}_{W} P(W|O) = \operatorname*{argmax}_{W} P(A|O)P(O|Q)P(Q|L)P(L|W)P(W)$$

| $P(A|O)$ | $P(O|Q)$ | $P(Q|L)$ | $P(L|W)$ | $P(W)$ |
|---|---|---|---|---|
| Feature Extraction | Frame Classification | Sequence Model | Lexicon Model | Language Model |

t ah m aa t ow

| $A$ | $O$ | $Q$ | $L$ | $W$ | |
|---|---|---|---|---|---|
| Speech Audio | Feature Frames | Sequence States | Phonemes | Words | Sentence |

http://slazebni.cs.illinois.edu/spring17/lec26_audio.pdf

# Before Deep Learning: Hidden Markov Models (HMM)

The Markov chain whose state sequence is unknown.

http://hts.sp.nitech.ac.jp/archives/2.3/HTS_Slides.zip

# Before Deep Learning: Gaussian Mixture Models (GMM)

Output probability is modeled by Gaussian mixture models.

http://hts.sp.nitech.ac.jp/archives/2.3/HTS_Slides.zip

# Before Deep Learning: N-gram Language Model

Given a sequence of N-1 words, an N-gram model predicts the most probable word that might follow this sequence.

## This is Big Data AI Book

| | | | | | |
|---|---|---|---|---|---|
| *Uni-Gram* | This | Is | Big | Data | AI | Book |

| | | | | |
|---|---|---|---|---|
| *Bi-Gram* | This is | Is Big | Big Data | Data AI | AI Book |

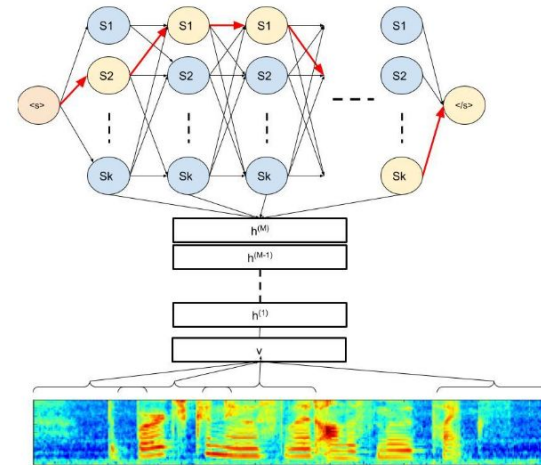| | | | |
|---|---|---|---|
| *Tri-Gram* | This is Big | Is Big Data | Big Data AI | Data AI Book |

https://devopedia.org/images/article/219/7356.1569499094.png

# Acoustic Modeling: GMM–HMM to DNN–HMM

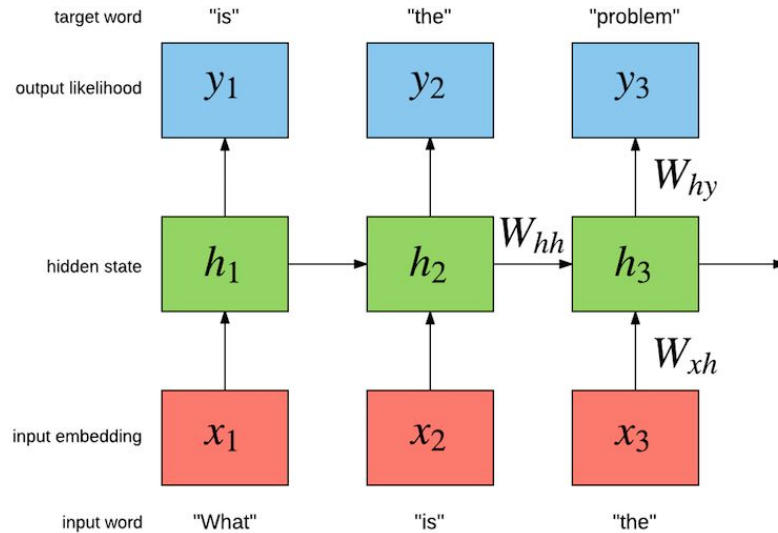Deep Neural Network (DNN) replaced Gaussian Mixture Model (GMM) for audio frame classification.
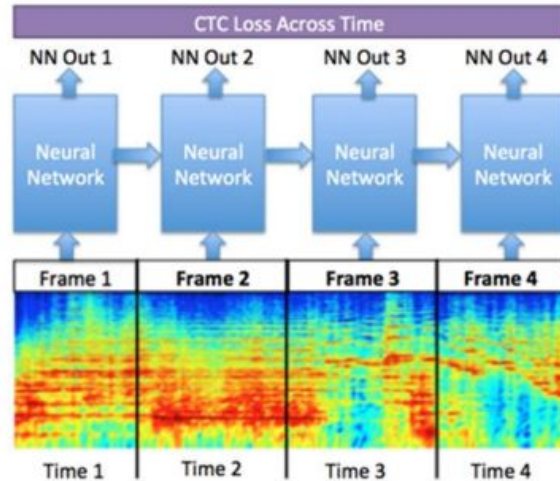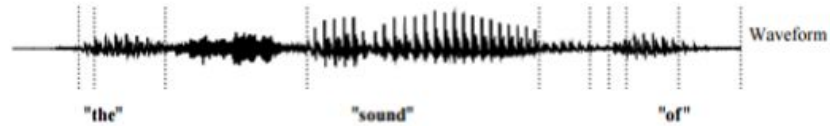
# Language Modeling: N-gram to NNLM

Neural network language model (NNLM) replaced N-gram for language modeling.

http://torch.ch/blog/2016/07/25/nce.html

# Sequence Modeling: HMM to CTC

Connectionist Temporal Classification (CTC) is a method for labeling unsegmented data sequences.
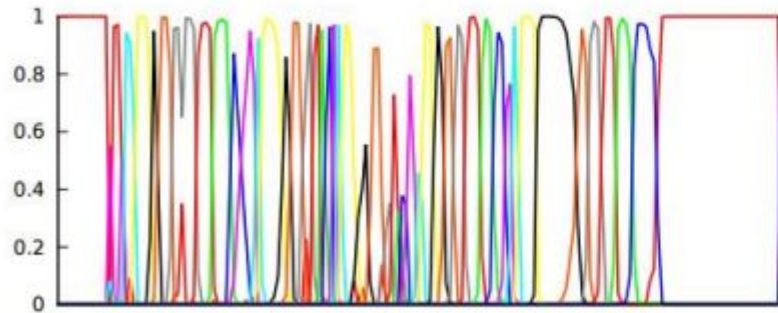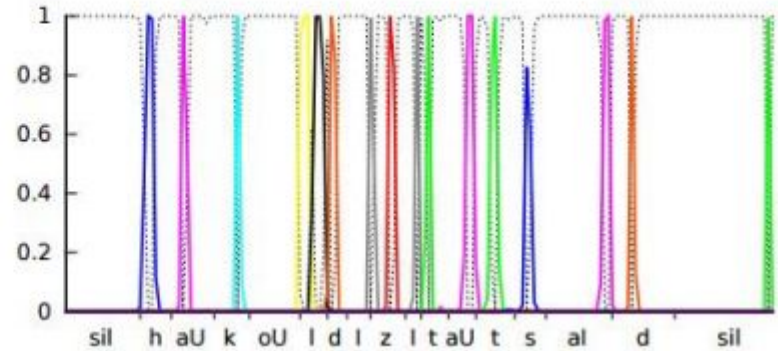
ftp://ftp.idsia.ch/pub/juergen/icml2006.pdf

# Sequence Modeling: HMM to CTC

Output probability of LSTM–HMM vs. LSTM–CTC.

CTC has spiky predictions, more discriminable between states than HMM.

Sak, Haşim, et al. "Learning acoustic frame labeling for speech recognition with recurrent neural networks." ICASSP, 2015.
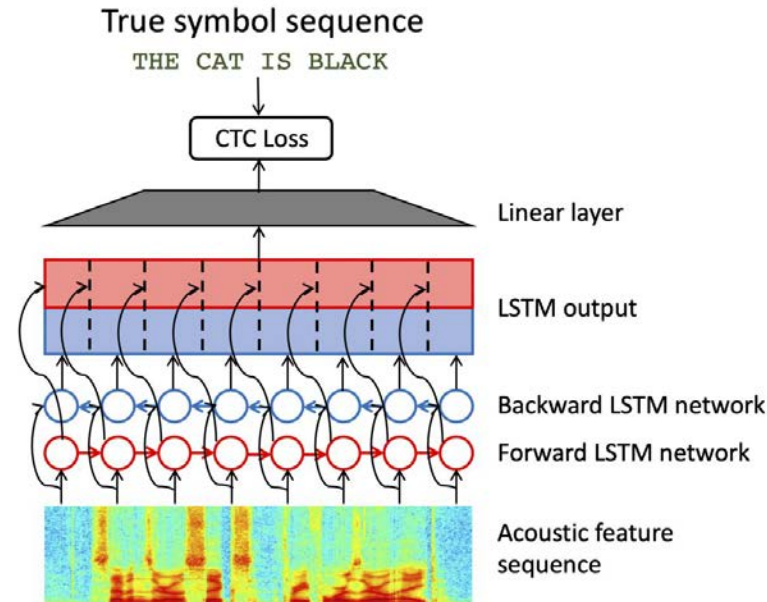
# Hybrid ASR to End-to-end ASR

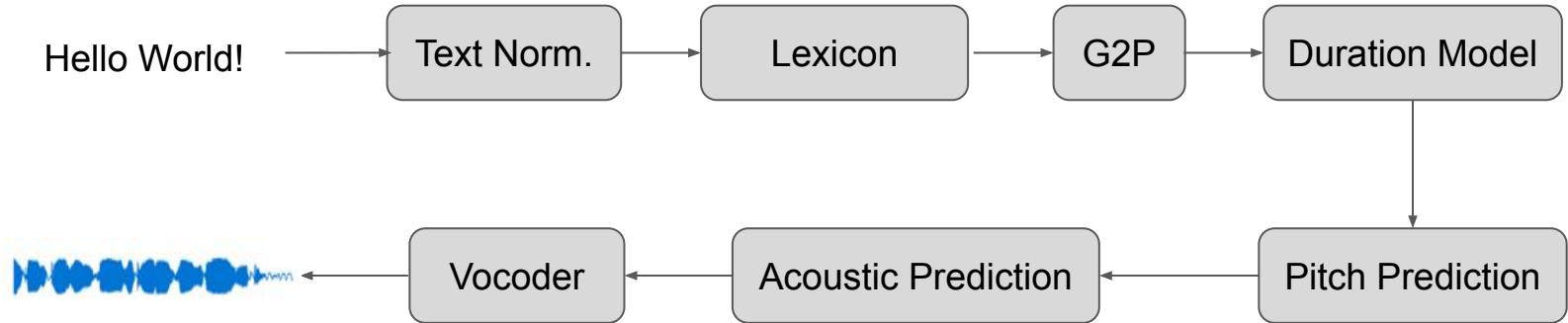**Hybrid ASR**: separated acoustic model, language model, and pronunciation model.

**End-to-end ASR**: a neural network to tackle audio frame sequence to word sequence mapping as a sequence-to-sequence learning problem.



https://www.ibm.com/blogs/research/2019/10/end-to-end-speech-recognition/

# Text-to-speech (TTS)

**Traditional Paradigms**:

- Concatenative synthesis
- Parametric synthesis

Hello World! → Text Norm. → Lexicon → G2P → Duration Model → Pitch Prediction → Acoustic Prediction → Vocoder →

# Sequence-to-sequence Modeling

- **Speech Recognition (continuous -> discrete):**

 → Hello World!

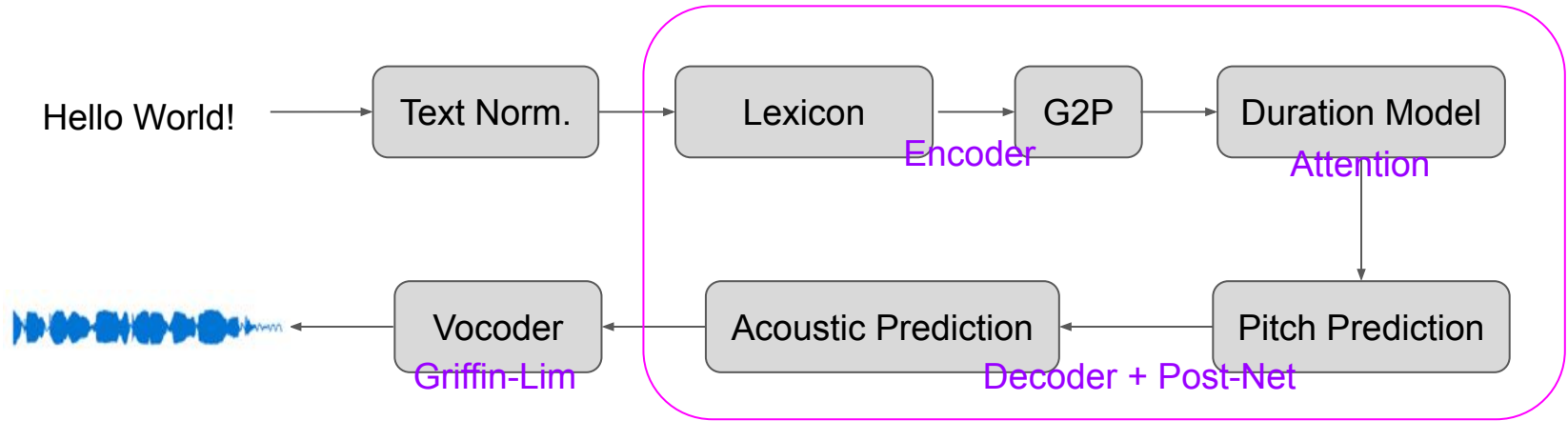- **Speech Synthesis (discrete -> continuous):**

Hello World! → 

- **Machine Translation (discrete -> discrete):**
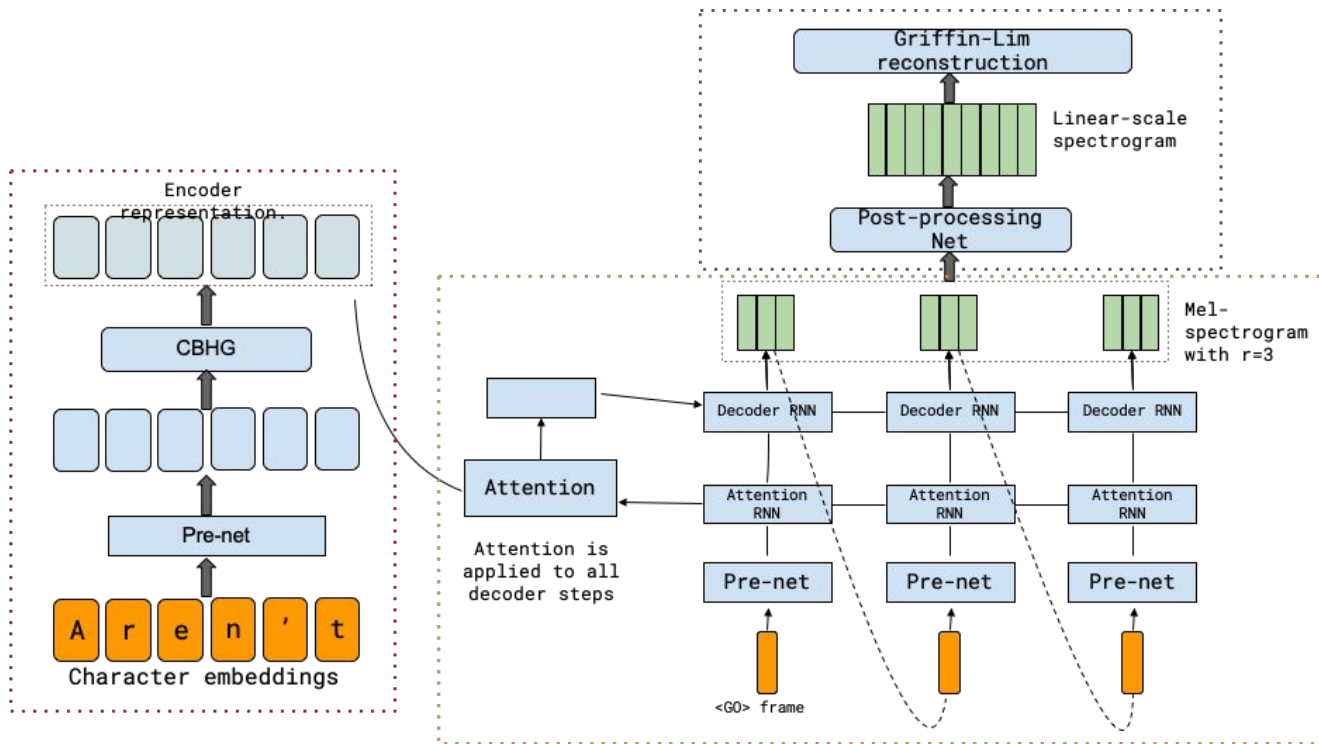
Hello World! → Hola Mundo!

# End-to-end TTS

A neural network to tackle word sequence to audio frame sequence mapping as a sequence-to-sequence learning problem.

ASAPP

# Tacotron: Google's End-to-end TTS architecture

https://arxiv.org/pdf/1703.10135.pdf

# Common NLP Tasks

NLP enables computers to perform a wide range of natural language related tasks at all levels, ranging from parsing and part–of–speech (POS) tagging, to machine translation and dialogue systems.

| Sentence Parsing | Word Tagging | Text Classification | Text Pair Matching | Text Generation |
|---|---|---|---|---|
| • Constituency parsing<br>• Semantic parsing<br>• Dependency parsing | • Word segmentation<br>• Shallow syntax-chunking<br>• Named entity recognition<br>• Part-of-speech tagging<br>• Semantic role labeling<br>• Word sense disambiguation | • Sentiment analysis<br>• Text classification<br>• Temporal processing<br>• Coreference resolution | • Semantic textual similarity<br>• Natural language inference<br>• Relation prediction | • Language modeling<br>• Machine translation<br>• Simplification<br>• Summarization<br>• Dialogue<br>• Question answering |

**ASAPP**

https://mobidev.biz/blog/natural-language-processing-nlp-use-cases-business

# An Example of NLP Linguistic Structure Analysis

**Part of speech:**

| NNP | NNP | RB | VBD | IN | NNP | NNP | , | CC | PRP | VBZ | RB | VBG | PRP | IN | PRP | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mrs. | Clinton | previously | worked | for | Mr. | Obama | , | but | she | is | now | distancing | herself | from | him | . |

**Named entity recognition:**

Mrs. Clinton [Person] previously [Date] worked for Mr. Obama [Person], but she is now [Date] distancing herself from him.

**Co-reference:**

[Mention] Mrs. Clinton previously worked for Mr. Obama [Ment], but she [M] is now distancing herself [Mention] from him [M]. — Coref

**Basic dependencies:**

https://cs224d.stanford.edu/papers/advances.pdf
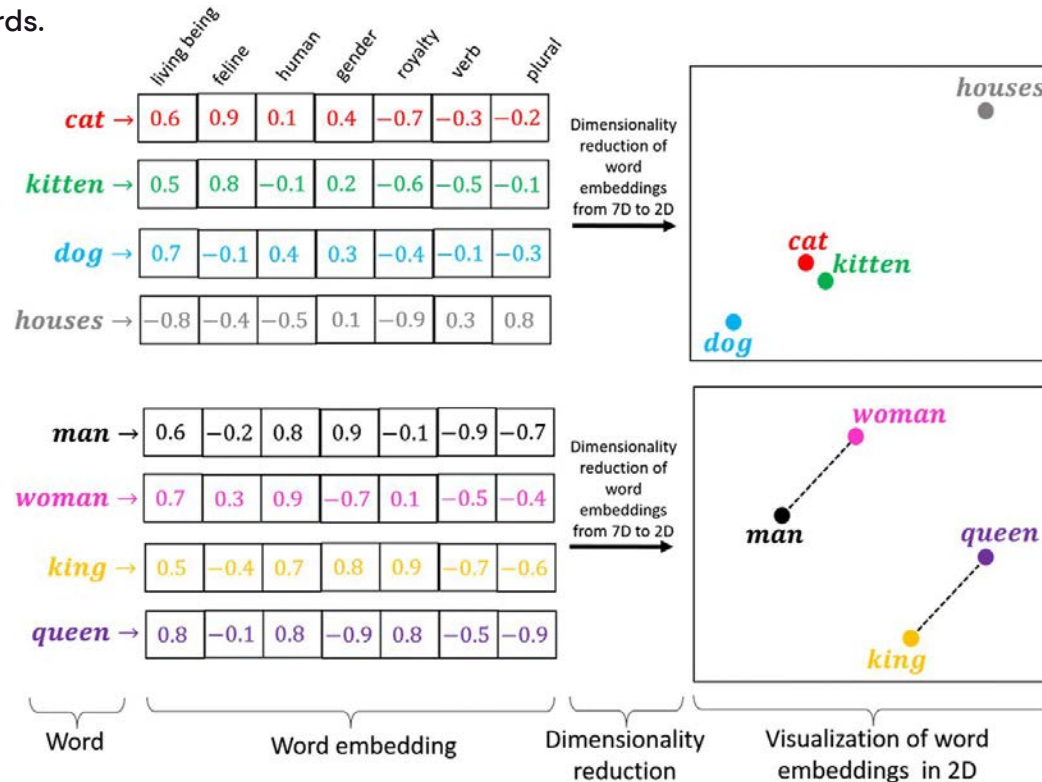
# An Example of NLP Question Answering



The first recorded travels by Europeans to China and back date from this time. The most famous traveler of the period was the Venetian Marco Polo, whose account of his trip to "Cambaluc," the capital of the Great Khan, and of life there astounded the people of Europe. The account of his travels, Il milione (or, The Million, known in English as the Travels of Marco Polo), appeared about the year 1299. Some argue over the accuracy of Marco Polo's accounts due to the lack of mentioning the Great Wall of China, tea houses, which would have been a prominent sight since Europeans had yet to adopt a tea culture, as well the practice of foot binding by the women in capital of the Great Khan. Some suggest that Marco Polo acquired much of his knowledge through contact with Persian traders since many of the places he named were in Persian.

How did some suspect that Polo learned about China instead of by actually visiting it?

**Answer:** through contact with Persian traders

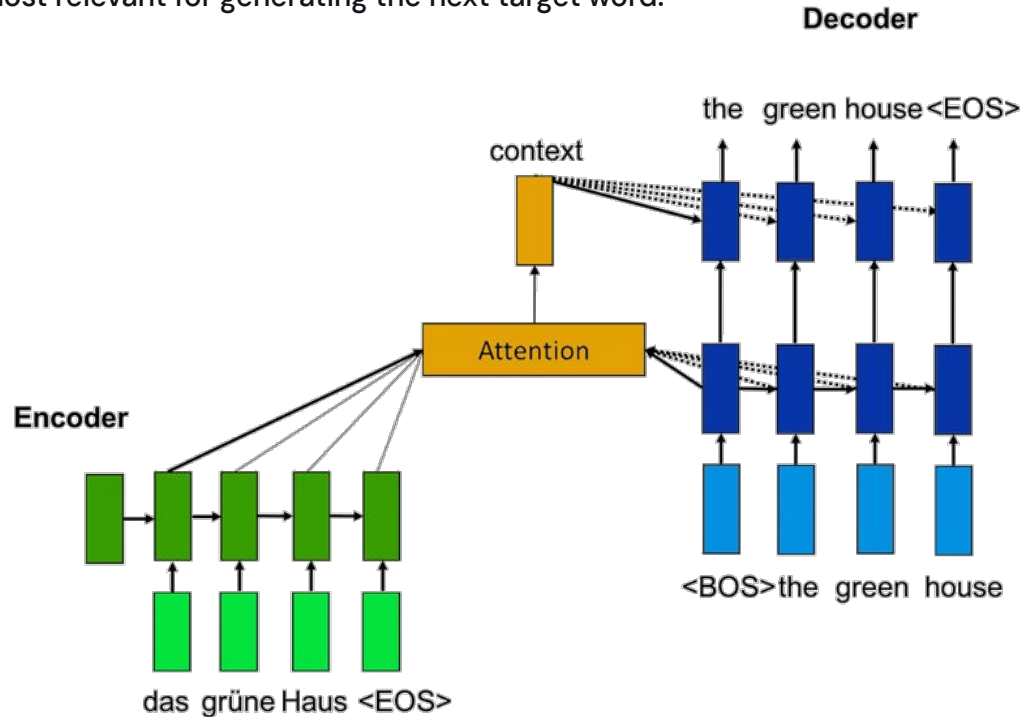https://rajpurkar.github.io/SQuAD-explorer/

# Word Embeddings

As the first data processing layer in a deep learning model, distributional vectors has advantage to capture similarity between words.

https://medium.com/@hari4om/word-embedding-d816f643140
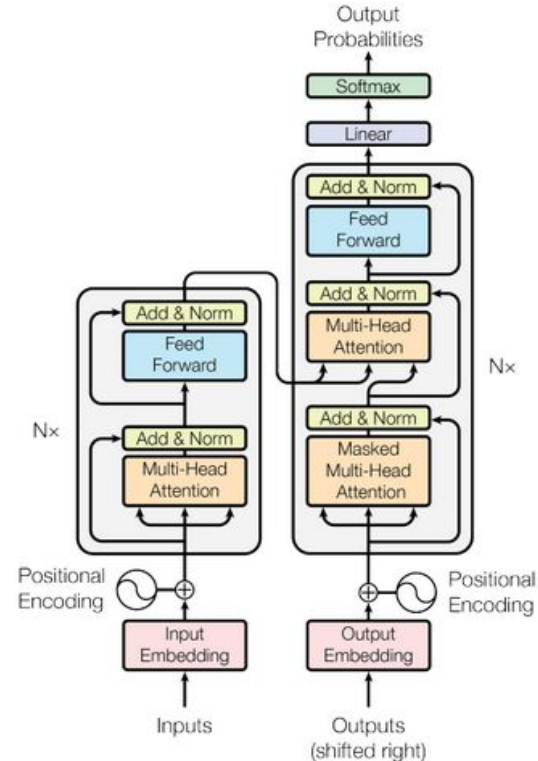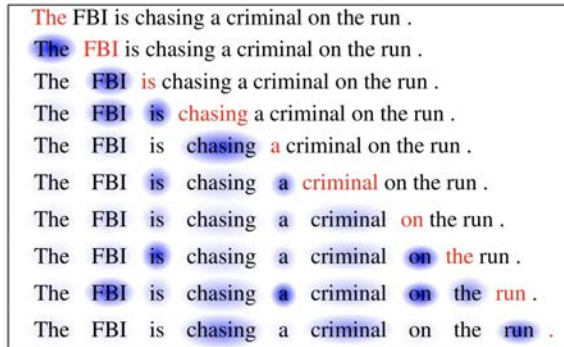
# Attention Mechanism

In a classic Neural Machine Translation model, using the attention mechanism, the decoder can decide which words are most relevant for generating the next target word.

https://aws.amazon.com/blogs/machine-learning/train-neural-machine-translation-models-with-sockeye/

# Self-attention & Transformer Architecture

**Self-attention**: allows the inputs to interact with each other and find out who they should pay more attention to.

**The Transformer:** first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution.
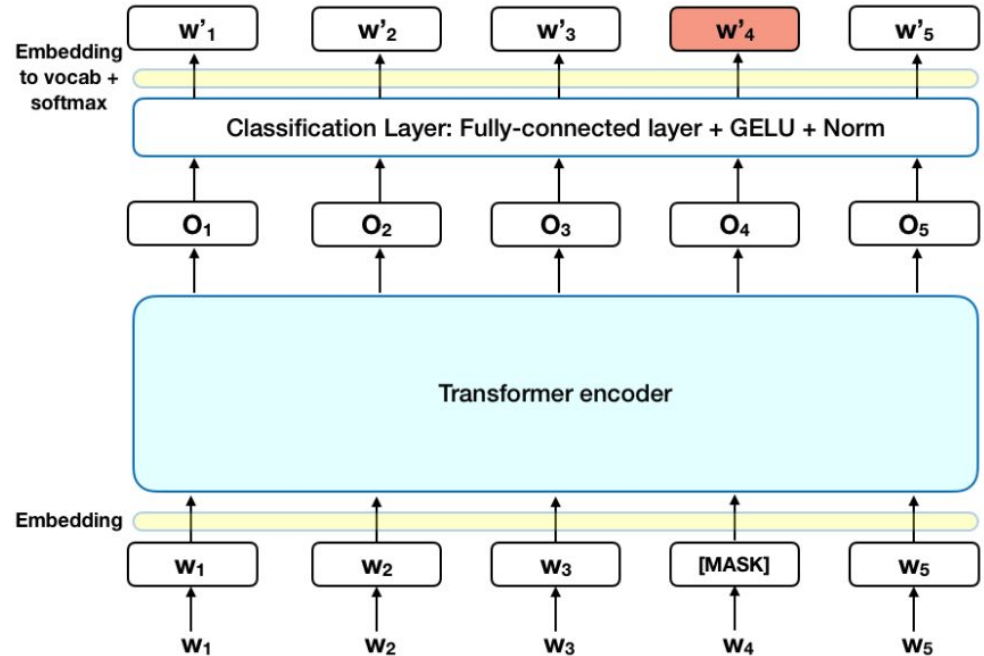




Vaswani et al. "Attention Is All You Need".
https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

# BERT: State-of-the-art Language Model for NLP

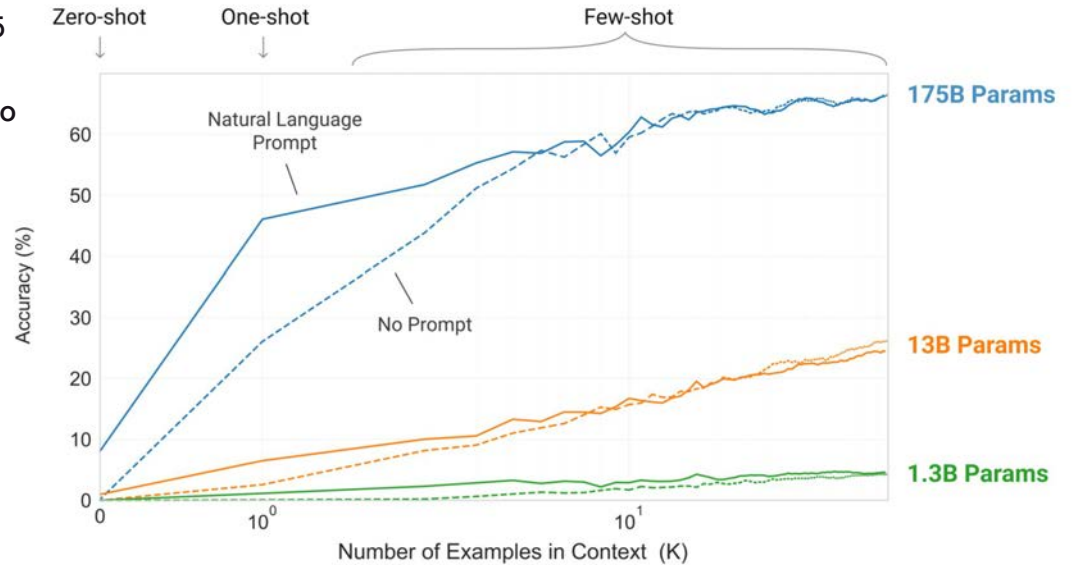**BERT**: Bidirectional Encoder Representations from Transformers, from Google.

Trained using Masked LM technique, BERT applied the bidirectional training of Transformer to language modelling and achieved state-of-the-art performance for many NLP tasks.

Devlin et al. "BERT: Pre-training of deep bidirectional transformers for language understanding."

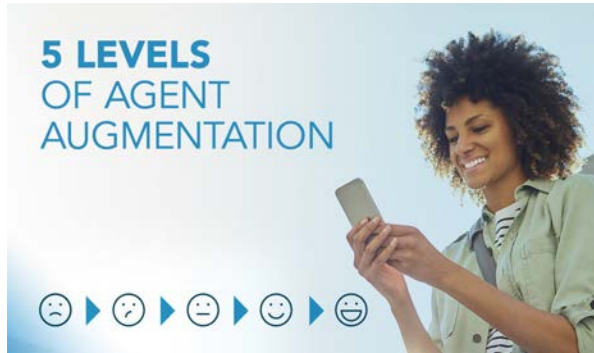# GPT-3: Predicts Anything that has a Language Structure

**GPT-3**: Generative Pre-trained Transformer 3, from OpenAI.

The largest language model ever trained, with 175 billion parameters. The quality of the text generated by GPT-3 is so high that it is difficult to distinguish from that written by a human.
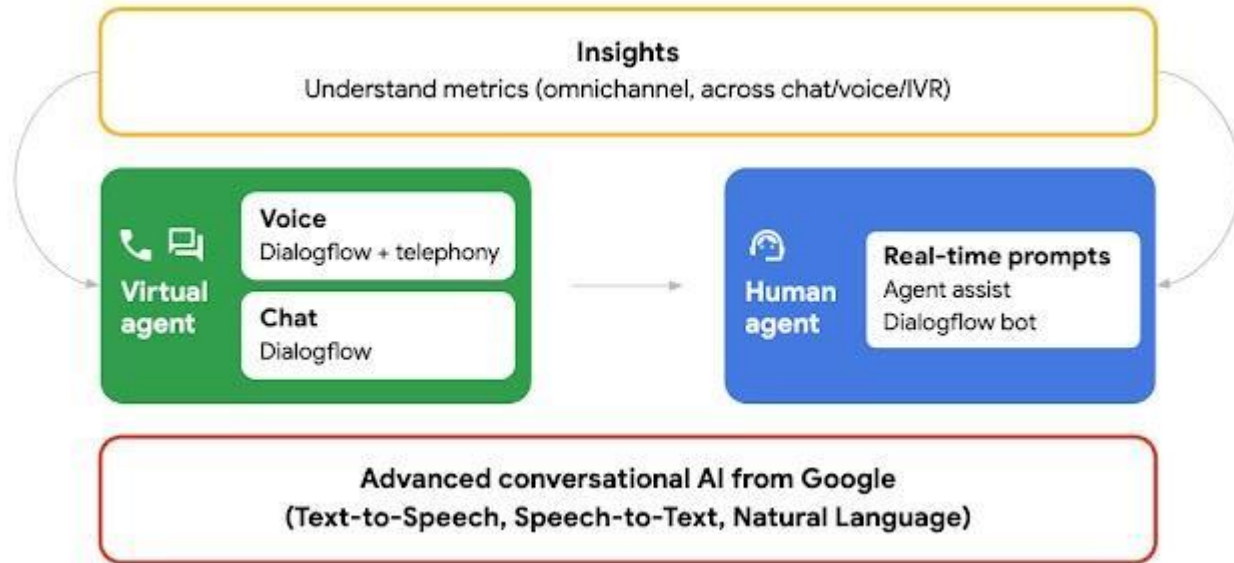
Brown et al. "Language Models are Few-Shot Learners".

# Language Technology Revolution is transforming Contact Center

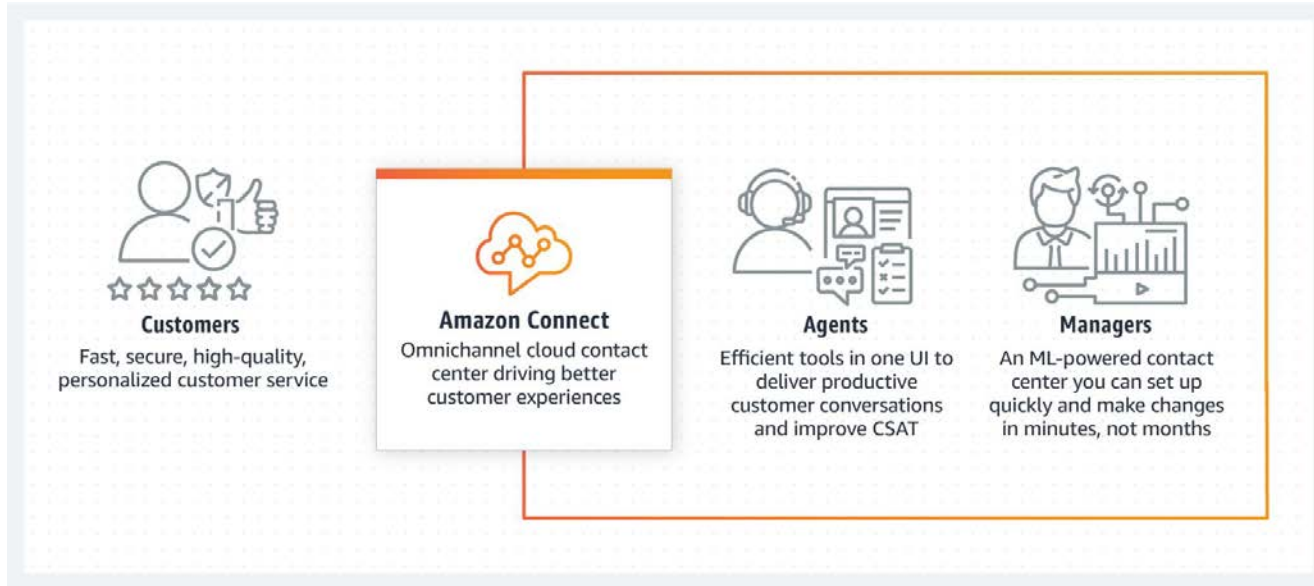Amid the COVID–19 crisis, the global market for Call Centers estimated at $339.4 Billion in the year 2020.

# Google Contact Center AI

Improve customer service with AI that understands, interacts, and talks.

https://cloud.google.com/solutions/contact-center

# Amazon Connect

An omnichannel cloud contact center that helps you provide superior customer service at a lower cost.

https://aws.amazon.com/connect/

# Future Directions & Opportunities

"*Any sufficiently advanced technology is indistinguishable from magic.*"
-- Arthur C. Clarke

### Joint Modeling of Speech & NLP

The *Status Quo* of separated Speech modeling and NLP modeling constraints language technology performance. Technology maturity of end-to-end deep learning architectures is opening up opportunities for a holistic language processing.

### Multimodal Communication

With contact center adopting new communication technologies such as WebRTC, a unified multimodal interaction platform across voice, text, and video starts emerging.

**ASAPP**

# Thank you!

ASAPP