

Generative AI for Dialog Systems
Text-to-Speech and Large Language Models



Themos Stafylakis
Head of Machine Learning and Voice Biometrics

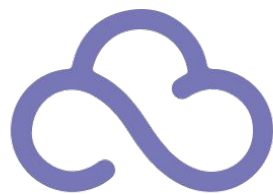
Outline

- Omilia Conversational AI
 - Omilia Cloud Platform (OCP)
 - Task-oriented Dialog Systems
 - ASR, NLU, Voice Biometrics, Dialog Manager, NLG & TTS
- Generative AI for Dialog Systems
 - Text-to-speech
 - From WaveNet to VALL-E
 - Large Language Models
 - Use-cases for task-oriented dialog systems

About Omilia

- Founded in 2002 by Dimitris Vassos (CEO) and Pelias Ioannidis (CFO).
- It creates automated dialog systems for customer care via voice, text, social media, a.o. communication channels.
- About 340 people mainly in Greece (Athens and Chania), Ukraine (Kiev), Cyprus (Limassol), Azerbaijan, Czech Republic, and USA.
- Omilia's clients are some of the largest financial institutions, telecoms, healthcare and insurance organizations, retail, e-commerce, restaurants, a.o. in USA, Canada, EU, Ukraine, a.o.
- During the last few years, Omilia is focusing on its cloud solution, Omilia Cloud Platform (OCP).





Omilia Cloud Platform

What is OCP[®]?

OCP[®] is a Conversational AI platform for customer care that

- requires no-code,
- offers enterprise-scale customer service automation,
- works across voice and text channels,
- offers hyper-tuned out-of-the-box industry models,
- and ready-to-go integration with most leading CCaaS platforms.



Architecture



Speech to Text

Translating the spoken language to accurate text



Natural Language Understanding

Understanding the semantics of the captured text (intent, entities, dialog acts, relations)



Voice Biometrics

Validating a user's identity through AI, using their voice characteristics.



Dynamic Dialog Management

Associating an intent with the desired outcome, planning dialog and tracking context



Pre-built Skills

Billions of real interactions, distilled into hyper-tuned ready to use models.



Text to Speech

Generating dynamic speech that sounds human-like

Pre-built Tasks, Domains, and Industry Modules

Omilia's pre-packaged solutions come with a fully preloaded intelligent Virtual Assistant providing out-of-the-box recognition and understanding of all key concepts for a specific domain and language.

These pre-build packs come with Concept Annotation Dictionaries, Rules, and Intents

Industry Skills



Banking



Healthcare



Telecoms



Retail



Travel



Insurance



Utilities



E-commerce



Banking Skills



Healthcare Skills



Telecom Skills

Ret...

+400 Intents Ready-to-Go

- Account_Balance
- List_Transaction
- Bill_Payments
- Activate_Card
- Loyalty_Rewards
- Transfer_Funds
- Branch_Locator

Virtual Agent Flows Ready-to-Go

I want to make a payment on my credit card

Allright. I see that you have 2 bank accounts registered for payment. Which one would you like to pay from?

The one ending in 2571

Ok. And how much would you like to pay?

My full balance

And when should I schedule the payment for?

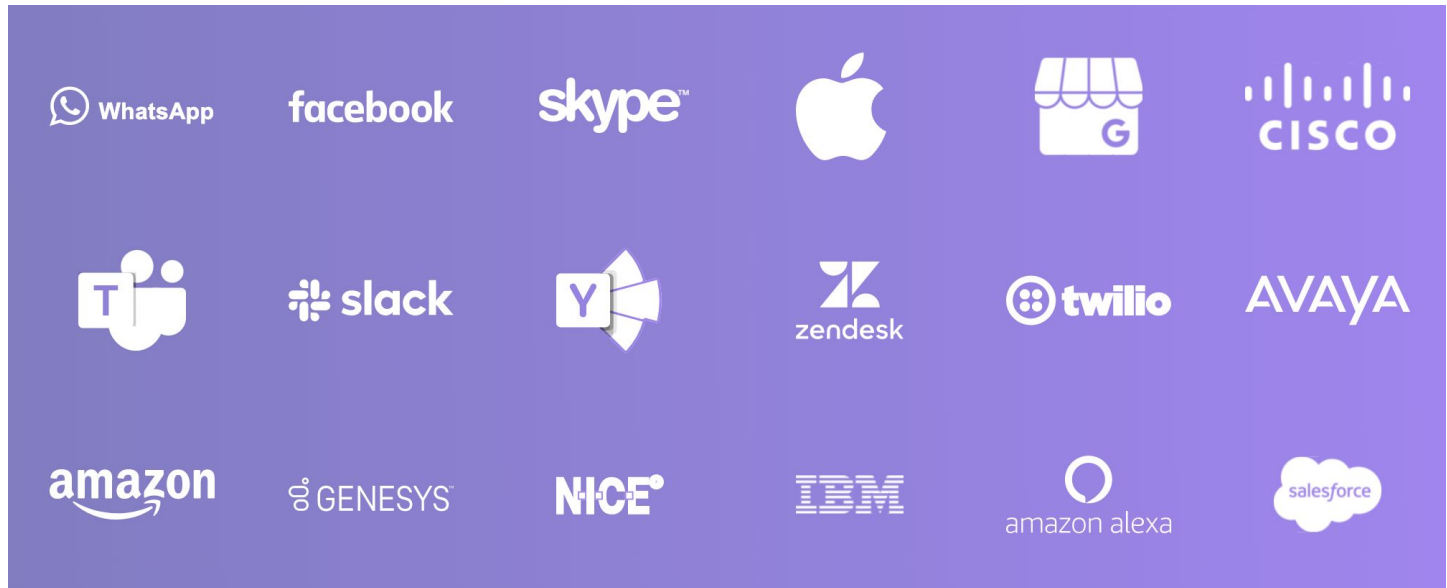
Today is fine.

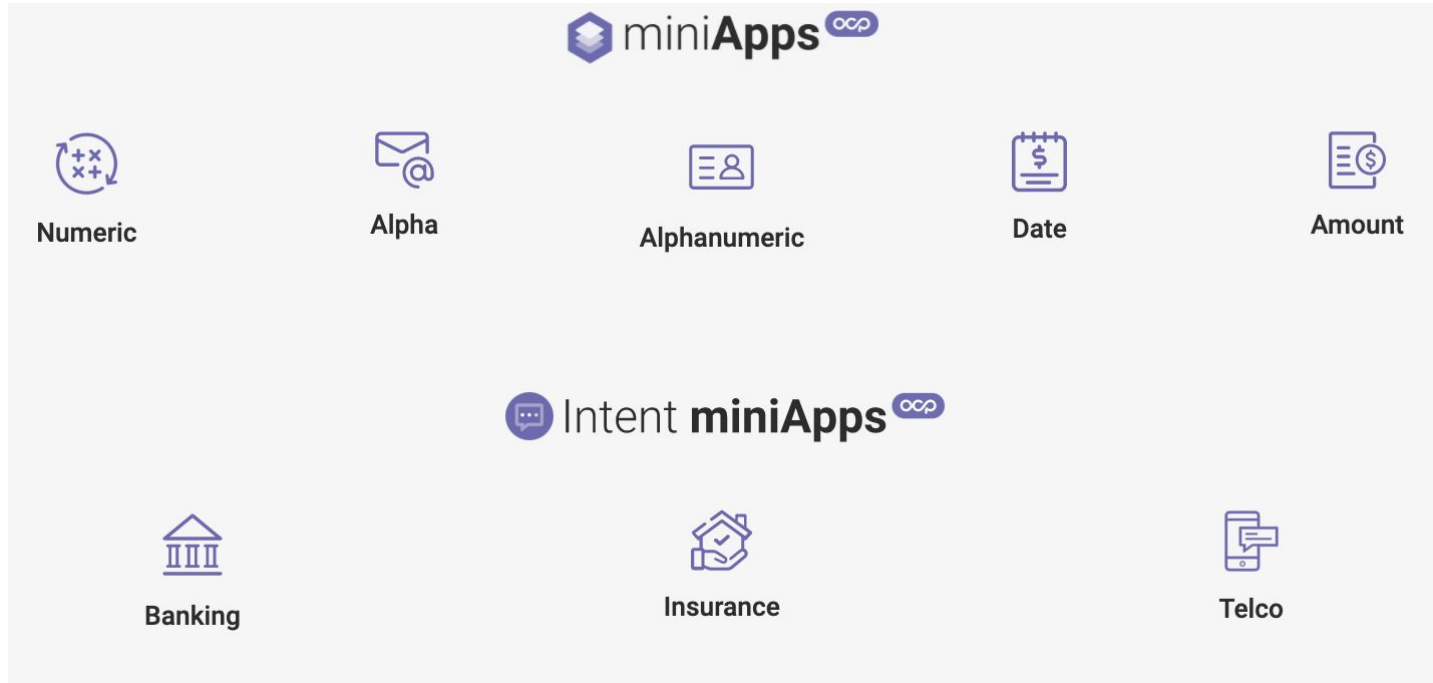
Let's recap: You've chosen to pay the total balance of \$XX with your BANKNAME TYPE account ending in XXXX, effective today/on #Date. Correct?

That is correct

Great! Your payment has been accepted and will be processed today. Your confirmation number is 123678. Is there anything else I can help you with?

Unlimited channel integration





Omilia Intent MiniApps[®]

The screenshot displays the Omilia NLU Models interface for a model named 'Petshop_Intent_Entity'. The interface is divided into a left sidebar and a main content area. The sidebar contains navigation options: 'NLU Models', 'Deployments', and a vertical menu of icons. The main content area shows the model's configuration, including its name, language (English (US)), and a 'Ready' status. Below this, there are tabs for 'Intents', 'Train', 'Evaluate', 'Upload Resources', 'Deployments', and 'Settings'. The 'Intents' tab is active, showing a list of 13 intents and a '+ Create' button. The 'Utterances' tab is also visible, showing 30 utterances. The 'Intents' list includes: Data_Privacy, My_Account, Order-Cancellation, Order-Details, Order-Return, Partner_Programme, Payment, Product_Details (highlighted in blue), Shop_Benefits, Technical_Issues, Card-Activation^{RB}, Card-Lost_Stolen^{RB}, and Card^{RB}. The 'Utterances' list includes: '+ Add utterance', 'Learn more about an item', 'Learn more about a product', 'Find some info about an item', 'Find some info about a product', 'Give me some info about an item', 'Give me some info about a product', 'Get some info about an item', 'Get some info about a product', 'Get some details about an item', 'Get some details about a product', 'I need to learn more about an item', 'I need to learn more about a product', 'I need to find some info about an item', 'I need to find some info about a product', 'I need to give me some info about an item', 'I need to give me some info about a product', and 'I need to get some info about an item'.

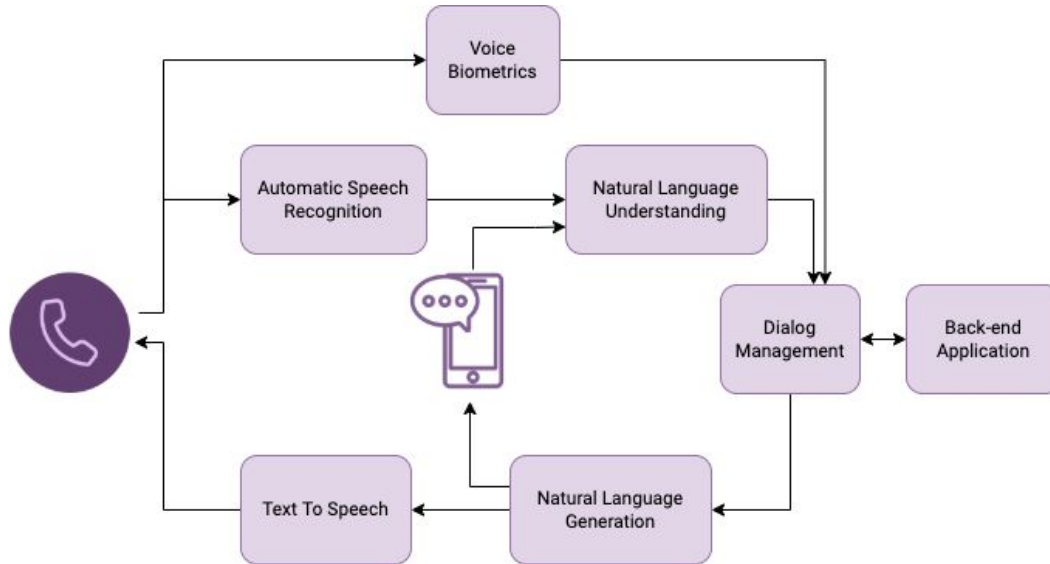
Omilia Intent MiniApps[®]

The screenshot displays the Omilia Intent MiniApps interface. On the left is a vertical navigation menu with icons for Home, Deployments, and various tool icons. The main content area shows the configuration for a model named 'Petshop_Intent_Entity' in English (US) with a custom ID: 238d673d-39e6-43af-a892-ca44480d985a. A 'Ready' status indicator is visible in the top right. Below the model name are tabs for Intents, Train, Evaluate, Upload Resources, Deployments, and Settings. The 'Train' tab is active, displaying a message: 'Your Model was Successfully Trained. More details below.' Below this message is a 'Train' button, a completion timestamp 'Finished on February 15th, 2022 14:55 UTC', and a 'Training Duration' indicator. A callout box explains: 'On Runtime, a Trained Model combines Machine Learning (using your Intents) and Pre-tuned xPort Package resources for Intent & Entity extraction.'

The screenshot displays the Orchestrator interface for a 'Petshop_rnd' application. On the left, a sidebar contains navigation options: miniApps (Manage), petshop_Tran... (Announcement), Technical_Lis... (Announcement), Transfer_Acc... (Announcement), Transfer_Ord... (Announcement), Flows, Dialog Control, and Transfer. The main workspace shows a flowchart starting with a 'petshop_inte... Universal' node leading to a 'Condition' node. This condition node branches into several paths based on intents: 'Intent = My_account', 'Intent = Data_Privacy', 'Intent = Technical_Issues', 'Intent = Order-Details', 'Intent = Payment', and 'Intent = Order-Return OR Intent = Order-Cancellation'. Each path leads to specific actions like 'Get_usernam...', 'Data_privacy Announcement', 'Technical_Lis... Announcement', 'petshop_ord... Alphanumeric', 'petshop_cust... Numeric', and 'petshop_Tran... Announcement'. The chat window on the right shows a conversation where a user asks for help, requests their order number, and asks to cancel an order. The virtual assistant responds with a welcome message, provides the order number 'ABC123', and offers to transfer the user to the Order Management Department.

Conversational AI

Voice-enabled task-oriented dialog systems



Automatic Speech Recognition

Omilia ASR Services

deepASR®

- Conversational Speech-to-Text (real-time)
- Transcription Service (off-line)

Deployment: Cloud, On-Premises

Domains: Financial Services (Banking, Insurance, Investment), Telecoms, Healthcare, Government, Retail/Wholesale, Manufacturing, Transportation (Travel & Hospitality), Utilities, E-Commerce

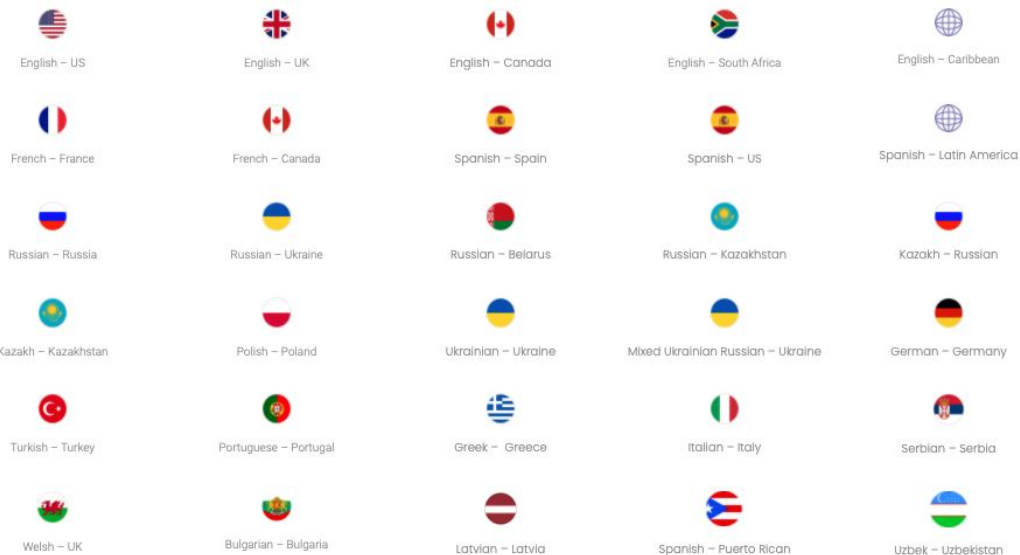
Features: Context-sensitive & dynamic recognition grammars, Rich transcriptions (Time stamps, Confidence Scores)

Customization: Acoustic Model, Language Model (word, sentence boosting), Lexicon (add new words)

Automatic Speech Recognition

ASR models in 30 languages and dialects

deepASR®



Voice Biometrics

Authenticate customers in the background as they speak

Enrollment phase

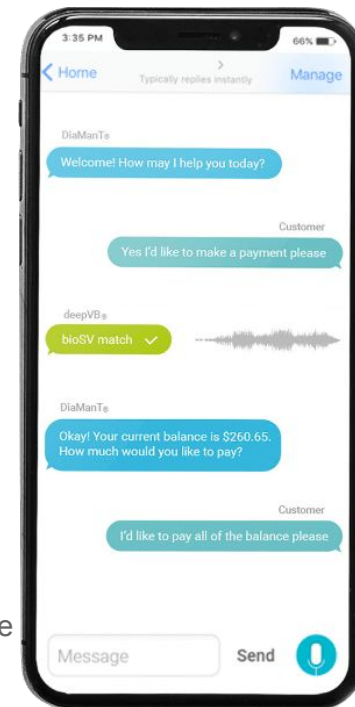
deepVB®

1. The callers are asked to enroll in the VB system.
2. If they opt-in, their voice sample is stored in the VB server and a **voiceprint** is created and stored in the database.

Note: The user may ask to unsubscribe from the VB system anytime.

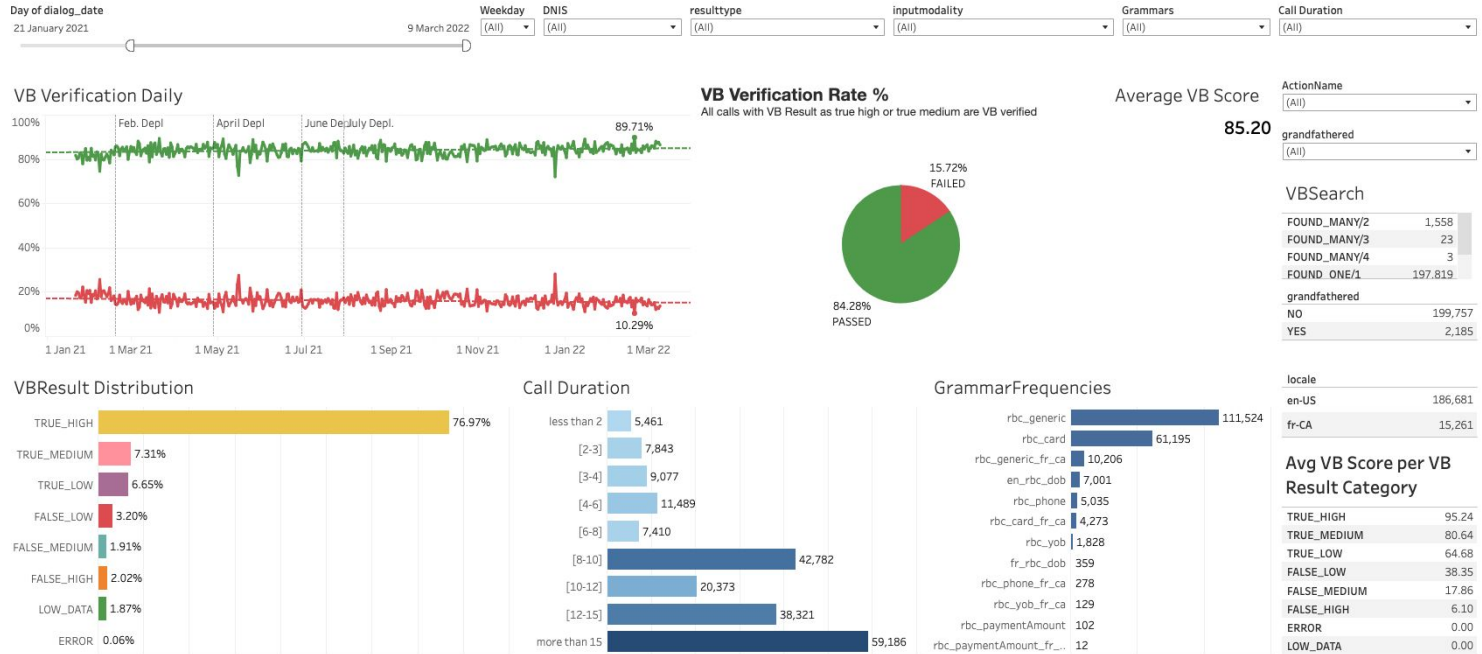
Verification phase

1. The VB system samples the speaker's voice and computes a voiceprint in real time.
2. It then compares the speaker's current voiceprint with the stored voiceprint and generates a "bioSV Confidence Score".
3. bioSV Confidence Score is used by the dialog manager to drive the conversation with the customer according to the business logic; deliver self-service, transfer to live agent, or ask additional traditional security questions.



Voice Biometrics

Tools for analysing VB results



Natural Language Understanding & Dialog Manager

Extracting meaning from utterances



DiaManT®



deepNLU®



xPert Packs®

- **Intents**
 - I want to access my bank account ⇨ [Access_AccountType](#)
 - I need online access ⇨ [Access_TelcoChannel](#)
 - I want to reactivate my accounts ⇨ [Activate_Account](#)
 - I call to cancel a warranty ⇨ [Cancel_Warranty](#)
- **Dialog Acts**
 - I need to speak to a customer service about logging into my account ⇨ [AgentRequest](#)
 - Balance again please ⇨ [RepeatRequest](#)
 - No, what I asked ⇨ [Rejection](#)
- **Entities**
 - Examples: [AccountType](#), [Airline](#), [amountToRedeem](#), [ComChannel](#), [DayOfMonth](#), [InsuranceType](#), [OfferPeriod](#), [SentimentDescriptor](#), [TransactionPeriod](#), [USZipCode](#)
- **Relations**
 - Examples in the domain of food ordering will follow

Natural Language Understanding & Dialog Manager



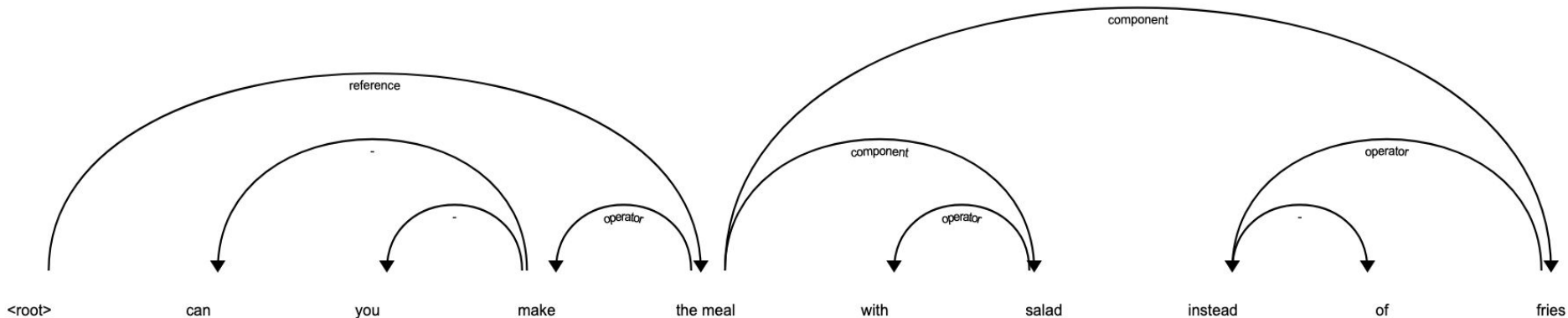
DiaManT®



deepNLU®

Extracting meaning from utterances

- **Relation extraction**
 - Examples in **food ordering**:
 - Can you make the meal with salad instead of fries?
- **Ontology**
 - Generic and/or restaurant-specific (e.g. derived from the menu)



Text-to-Speech

Transform text to speech

- English



- Ukrainian



- Russian

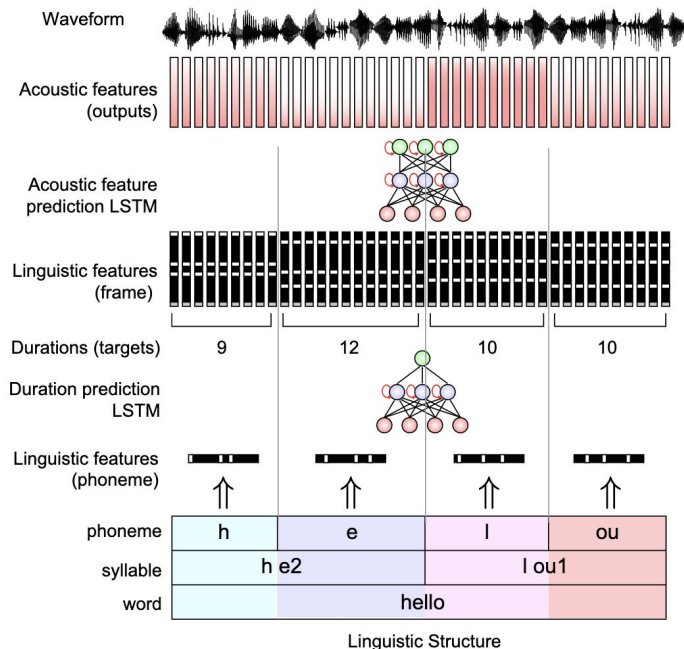


- **Neural TTS synthesis**
- **Conversational TTS**
 - The voices are geared to be used in conversation in real-world applications
- **Ability to control prosody based on punctuation**
 - Do you want to cancel your card because it is lost or stolen?
 - Do you want to cancel your card because it is lost, or stolen?
- **Ability to voice important information (dates, digits, addresses, etc.)**
- **Ability to control emphasis in order to highlight important information in a conversation**
 - E.g. The four digit number you need is **5432**.
- **French and Spanish** are also supported.

Generative AI for Dialog Systems

Text-to-Speech

Text-to-Speech - Statistical Parametric Speech Synthesis



H. Zen et al (Google AI). Fast, Compact, and High Quality LSTM-RNN-Based Statistical Parametric Speech Synthesizers for Mobile Devices, 2016

Text-to-Speech - WaveNet (2016)

WAVENET: A GENERATIVE MODEL FOR RAW AUDIO

Aäron van den Oord

Sander Dieleman

Heiga Zen[†]

Karen Simonyan

Oriol Vinyals

Alex Graves

Nal Kalchbrenner

Andrew Senior

Koray Kavukcuoglu

{avdnoord, sedielem, heigazen, simonyan, vinyals, graves, nalk, andrewsenior, korayk}@google.com
Google DeepMind, London, UK

[†] Google, London, UK

Text-to-Speech - WaveNet (2016)

Abstract

- A **deep neural network** for generating raw **audio waveforms**.
- It is fully probabilistic and **autoregressive**:
 - It models the predictive distribution for each audio sample conditioned on all previous ones:

$$p(\mathbf{x}) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1})$$

- It can be **efficiently trained** on data with tens of thousands of samples per second of audio.
- Significantly more natural sounding than the best parametric and concatenative systems.
- A single WaveNet can capture the characteristics of **many different speakers**:
 - It can switch between them by conditioning on the **speaker identity**.
- It can even be trained to **model music** or to do **phoneme recognition**.

Text-to-Speech - WaveNet (2016)

Dilated Convolutions

- Increase the receptive field while keeping the number of parameters fixed.
- Also used in ASR, e.g. time-delayed NNet (TDNN in Kaldi).
- Convolutions are much faster in training than recurrent models (e.g. LSTMs or GRUs).
- Why? Teacher forcing.
- However generation is still very slow, due to its autoregressive architecture (no gain over RNNs).

$$p(\mathbf{x}) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1})$$

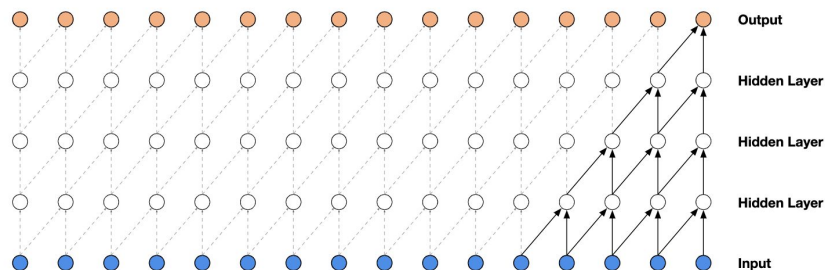


Figure 2: Visualization of a stack of causal convolutional layers.

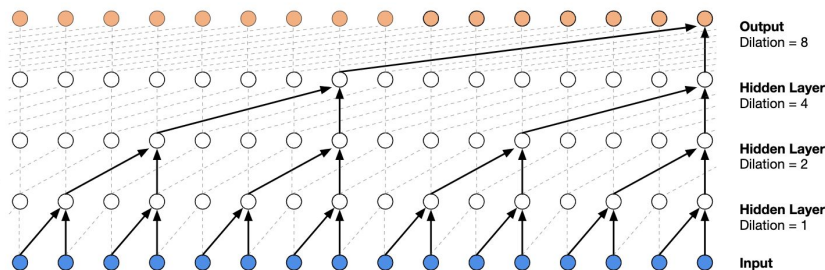
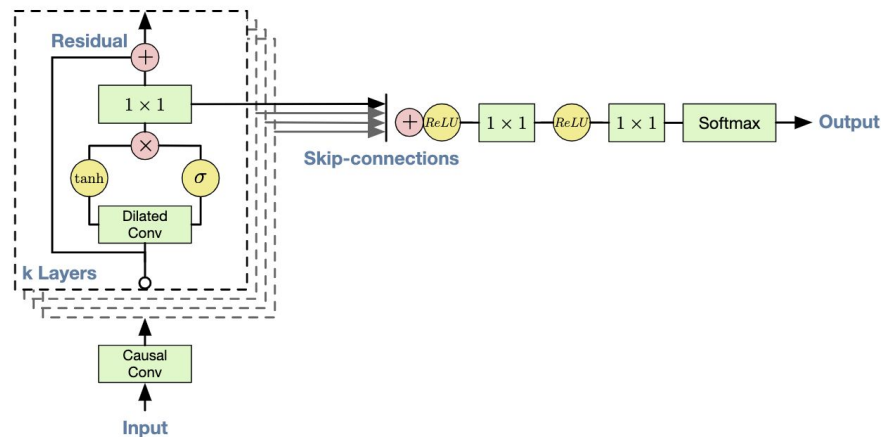


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

Text-to-Speech - WaveNet (2016)

Architecture



- Residual Connections: Each layer adds details to the reconstructed signal
- Gated Activation Units: $\mathbf{z} = \tanh(W_{f,k} * \mathbf{x}) \odot \sigma(W_{g,k} * \mathbf{x})$
- Softmax: Input and outputs are quantized to 256 levels/classes

Text-to-Speech - WaveNet (2016)

Conditional WaveNet

$$p(\mathbf{x} | \mathbf{h}) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}, \mathbf{h})$$

Condition Wavenet to mimic the voice of a speaker (\mathbf{h} is a speaker embeddings):

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k}^T \mathbf{h}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k}^T \mathbf{h})$$

or to follow a given text:

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k} * \mathbf{y}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k} * \mathbf{y})$$

where \mathbf{y} are the (time-synchronous) linguistic features of the desired text (a sequence indicating the phonetic class of each output sample, the prosody, a.o.).

Text-to-Speech - WaveNet (2016)

Speech samples	Subjective 5-scale MOS in naturalness	
	North American English	Mandarin Chinese
LSTM-RNN parametric	3.67 ± 0.098	3.79 ± 0.084
HMM-driven concatenative	3.86 ± 0.137	3.47 ± 0.108
WaveNet (L+F)	4.21 ± 0.081	4.08 ± 0.085
Natural (8-bit μ -law)	4.46 ± 0.067	4.25 ± 0.082
Natural (16-bit linear PCM)	4.55 ± 0.075	4.21 ± 0.071

Table 1: Subjective 5-scale mean opinion scores of speech samples from LSTM-RNN-based statistical parametric, HMM-driven unit selection concatenative, and proposed WaveNet-based speech synthesizers, 8-bit μ -law encoded natural speech, and 16-bit linear pulse-code modulation (PCM) natural speech. WaveNet improved the previous state of the art significantly, reducing the gap between natural speech and best previous model by more than 50%.

Text-to-Speech - WaveNet (2016)

WaveNet is very slow to be used in production, so why is it important?

- It was the first neural net to synthesize waveform (**features-to-waveform**).
 - Most TTS systems have a decoder (text-to-spectrogram) followed by a vocoder (spectrogram-to-waveform).
 - Prior to WaveNet the vocoder was typically **Griffin-Lim**.
- New **non-autoregressive models arrived**, overcoming the issue of slow inference:
 - **Parallel WaveNet** and **ClariNet** use a **pretrained WaveNet** for knowledge **distillation**.
 - 20 times faster than real time.
- **Offline generation** is still **useful**:
 - WaveNet was used to generate **Google Assistant** voices for US English and Japanese across all Google platforms.
- WaveNet can also be **conditioned on spectrogram**, i.e. it become a **vocoder**.
 - Many state-of-the-art TTS systems use a WaveNet-like architecture as a vocoder.
- It is **multispeaker**:
 - You may train a single WaveNet (with utterances of many speakers) and generate speech from any of them.

Text-to-Speech - Tacotron 2 (2017)

NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS

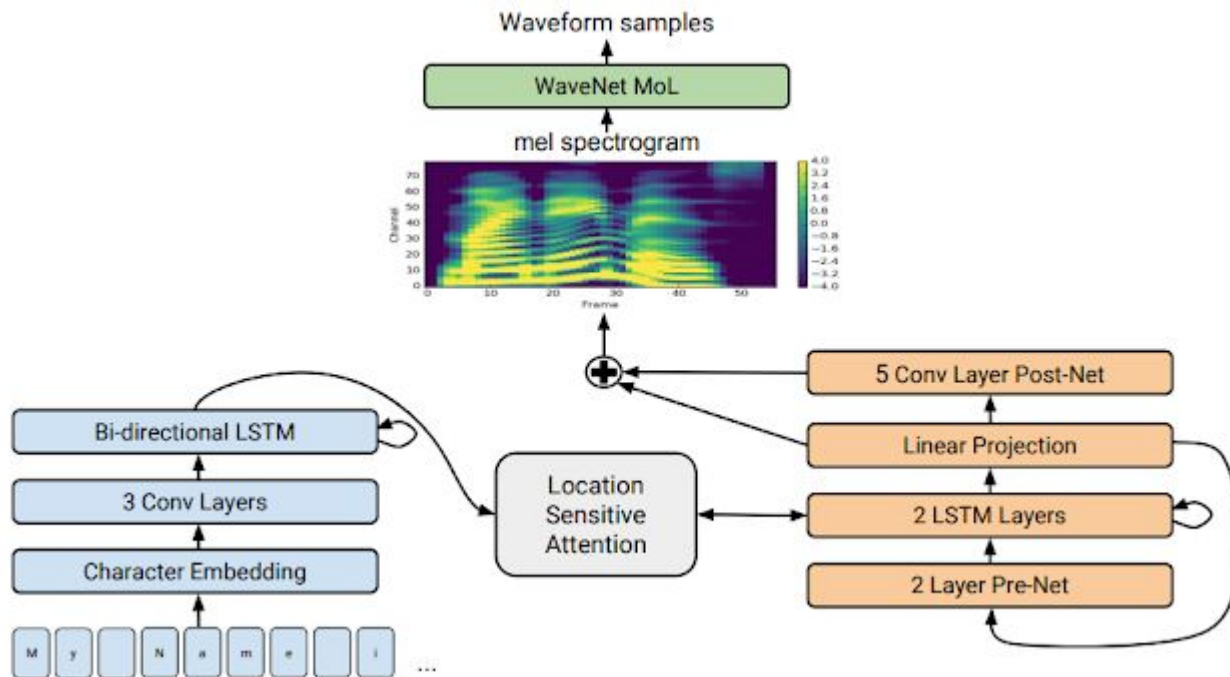
*Jonathan Shen¹, Ruoming Pang¹, Ron J. Weiss¹, Mike Schuster¹, Navdeep Jaitly¹, Zongheng Yang^{*2}, Zhifeng Chen¹, Yu Zhang¹, Yuxuan Wang¹, RJ Skerry-Ryan¹, Rif A. Saurous¹, Yannis Agiomyriannakis¹, and Yonghui Wu¹*

¹Google, Inc., ²University of California, Berkeley,
{jonathanasdf, rpang, yonghui}@google.com

Text-to-Speech - Tacotron 2 (2017)

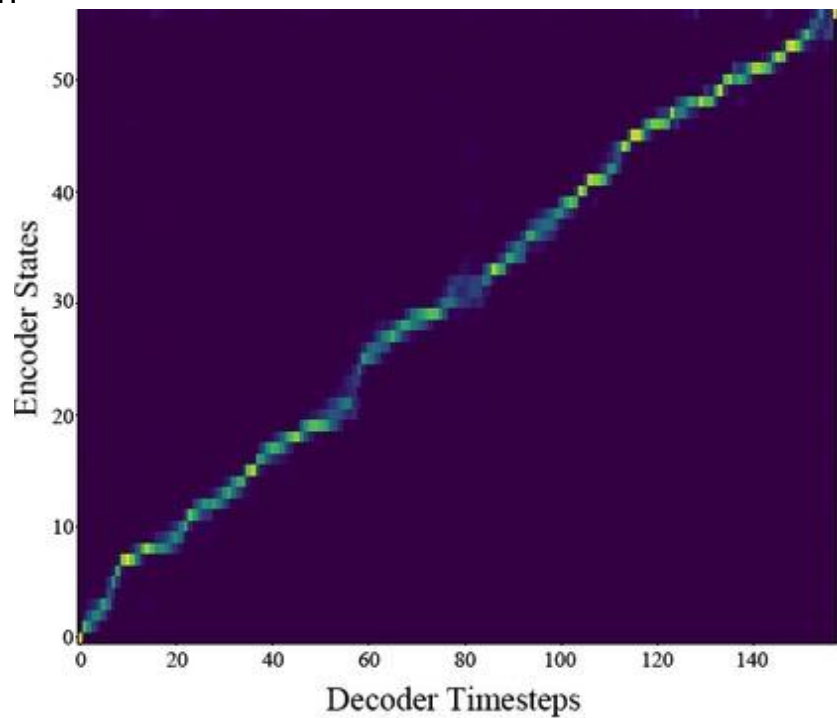
- It does not use complex linguistic and acoustic features as input.
- It comprises:
 - a recurrent **sequence-to-sequence** feature prediction network with attention which predicts a sequence of mel spectrogram frames from an input character sequence, and
 - a modified version of **WaveNet** which generates time-domain waveform samples conditioned on the predicted mel spectrogram frames.
- The alignment between the input character sequence and acoustic features is estimated by the **attention** mechanism.

Text-to-Speech - Tacotron 2 (2017)



Text-to-Speech - Tacotron 2 (2017)

Alignment via attention



Text-to-Speech - Tacotron 2 (2017)

- It works well on out-of-domain and complex words.
- It learns pronunciations based on phrase semantics.
- It is somewhat robust to spelling errors.
- It is sensitive to punctuation.
- It learns stress and intonation.
- Its prosody changes when turning a statement into a question.
- It is good at tongue twisters.

Samples:

<https://google.github.io/tacotron/publications/tacotron2/index.html>

Text-to-Speech - Zero-shot Multispeaker Tacotron 2 (2018)

Transfer Learning from Speaker Verification to Multispeaker Text-To-Speech Synthesis

Ye Jia* **Yu Zhang*** **Ron J. Weiss*** **Quan Wang** **Jonathan Shen** **Fei Ren**
Zhifeng Chen **Patrick Nguyen** **Ruoming Pang** **Ignacio Lopez Moreno** **Yonghui Wu**
Google Inc.
{jiaye,ngyuzh,ronw}@google.com

Text-to-Speech - Zero-shot Multispeaker Tacotron 2 (2018)

It can mimic the voice of speakers unseen during training

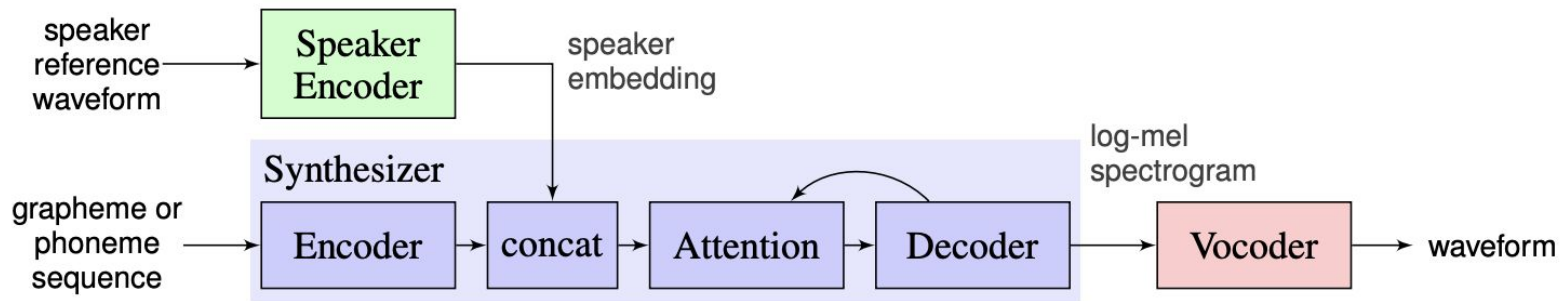


Figure 1: Model overview. Each of the three components are trained independently.

Text-to-Speech - Zero-shot Multispeaker Tacotron 2 (2018)

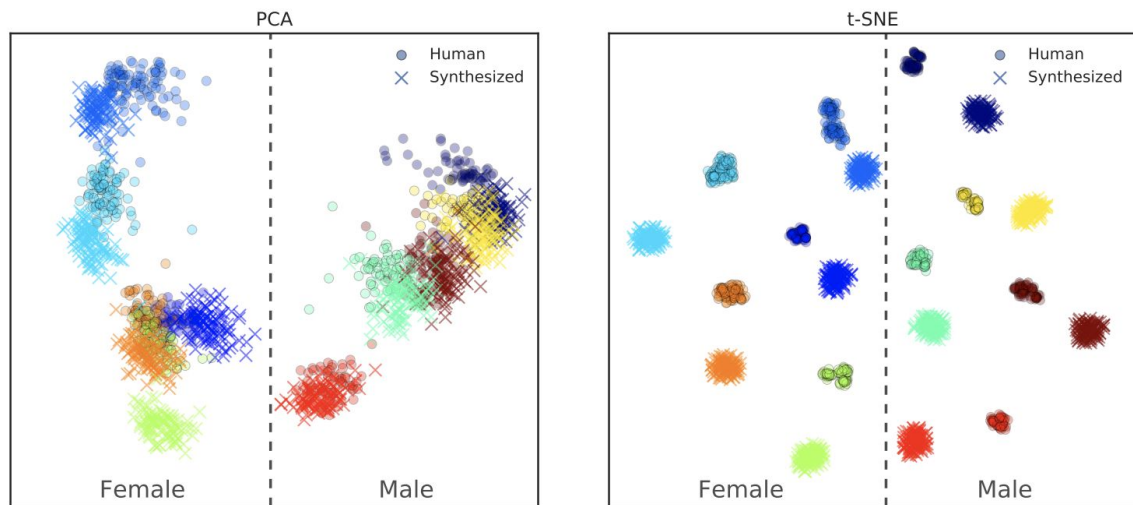


Figure 3: Visualization of speaker embeddings extracted from LibriSpeech utterances. Each color corresponds to a different speaker. Real and synthetic utterances appear nearby when they are from the same speaker, however real and synthetic utterances consistently form distinct clusters.

Samples: https://google.github.io/tacotron/publications/speaker_adaptation/

Text-to-Speech - FastSpeech 1&2 (2018)

FASTSPEECH 2: FAST AND HIGH-QUALITY END-TO-END TEXT TO SPEECH

Yi Ren^{1*}, Chenxu Hu^{1*}, Xu Tan², Tao Qin², Sheng Zhao³, Zhou Zhao^{1†}, Tie-Yan Liu²

¹Zhejiang University

{rayeren, chenxuhu, zhaozhou}@zju.edu.cn

²Microsoft Research Asia

{xuta, taoqin, tyliu}@microsoft.com

³Microsoft Azure Speech

Sheng.Zhao@microsoft.com

Text-to-Speech - FastSpeech 1&2 (2018)

- **FastSpeech 1**

- A **non-autoregressive** model trained using knowledge distillation from Tacotron 2.
- As a non-autoregressive model it is very fast in inference.
- However, **knowledge distillation** is slow in training (e.g. requires training a Tacotron 2).
- The duration extracted from the teacher model is not accurate enough.
- The target mel-spectrograms distilled from teacher model suffer from information loss due to data simplification.

- **FastSpeech 2**

- **Trained** directly with **ground-truth** targets.
- Introduces more variation information of speech (e.g., pitch, energy and more accurate duration) as conditional inputs.
- A very popular network for real-time production systems.

Text-to-Speech - FastSpeech 1&2 (2018)

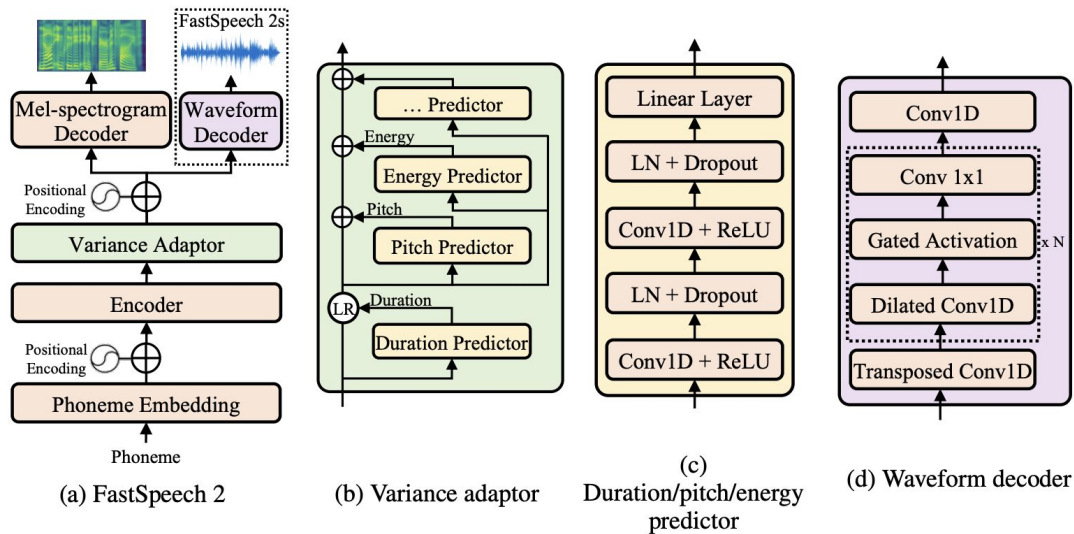
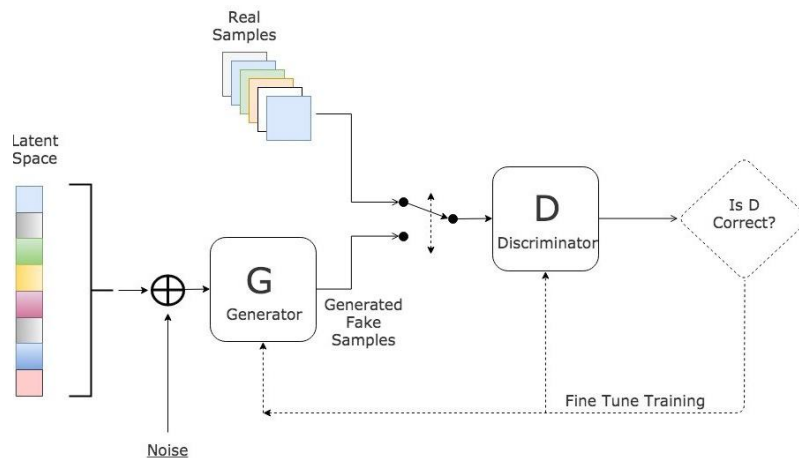


Figure 1: The overall architecture for FastSpeech 2 and 2s. LR in subfigure (b) denotes the length regulator proposed in FastSpeech. LN in subfigure (c) denotes layer normalization.

Samples: <https://speechresearch.github.io/fastspeech2/>

Vocoders using Generative Adversarial Nets (GANs)

- K. Kumar, et al. “MelGAN: Generative adversarial networks for conditional waveform synthesis,” in NeurIPS, 2019.
- M. Binkowski, et al. “High fidelity speech synthesis with adversarial networks,” in ICLR, 2020
- J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” in NeurIPS, 2020.



Generative Adversarial Nets

Text-to-Speech - VALL-E (2022)

Neural Codec Language Models are Zero-Shot Text to Speech Synthesizers

**Chengyi Wang* Sanyuan Chen* Yu Wu* Ziqiang Zhang Long Zhou Shujie Liu
Zhuo Chen Yanqing Liu Huaming Wang Jinyu Li Lei He Sheng Zhao Furu Wei**

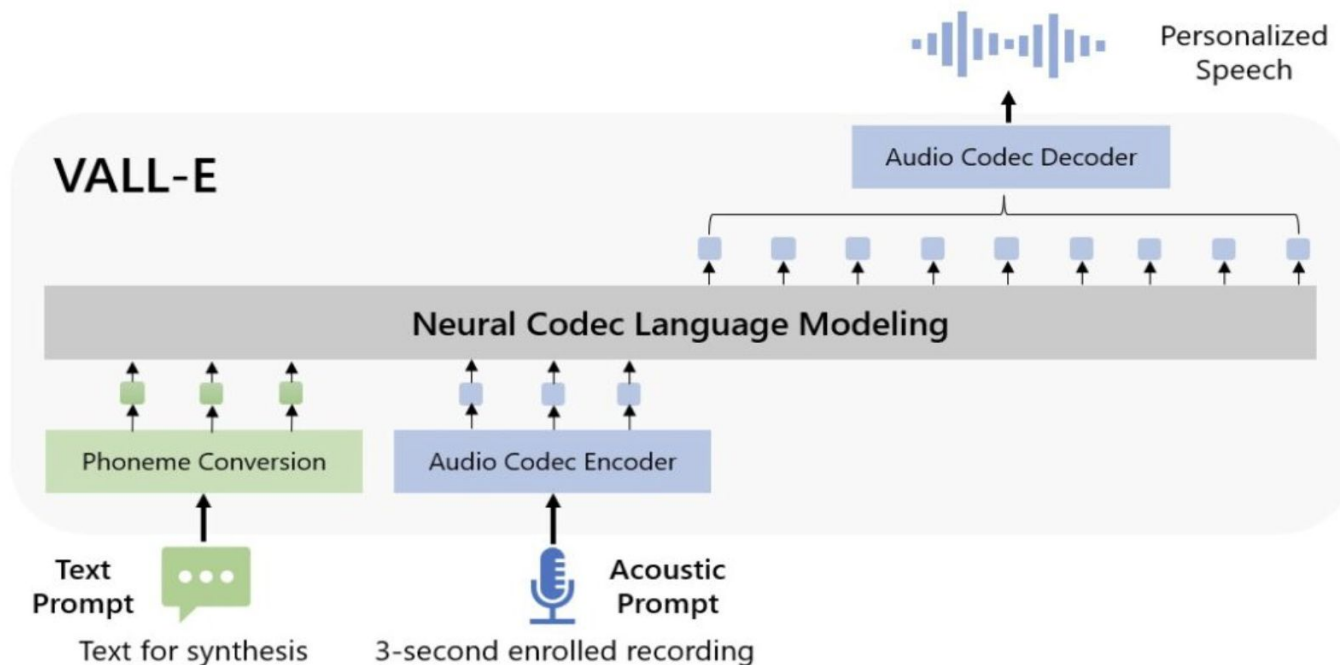
Microsoft

<https://github.com/microsoft/unilm>

Text-to-Speech - VALL-E (2022)

- VALL-E is a TTS framework by **Microsoft**.
- It treats TTS as a language model task and uses **audio codec** codes as an intermediate representation instead of mel spectrograms.
- The framework leverages a large amount of **semi-supervised data** to build a generalized TTS system in the speaker dimension, highlighting the importance of scaling up semi-supervised data for TTS.
- It uses a **transformer decoder** (instead of recurrent or convolutional).
- VALL-E can produce diverse outputs for the same input text while maintaining the **acoustic environment** and the **speaker's emotion** from the acoustic prompt.
- In a zero-shot scenario, VALL-E demonstrates natural speech synthesis with **high speaker similarity** compared to other systems.

Text-to-Speech - VALL-E (2022)



Samples: <https://www.microsoft.com/en-us/research/project/vall-e/>

Can we detect synthetic speech?

- Doable with Deep Neural Networks trained to classify real vs synthetic speech
- Fairly easily for known and old TTS systems
- Harder for new systems which are unseen during training



asvspoof.org

Generative AI for Dialog Systems

Large Language Models

Large Language Models

5

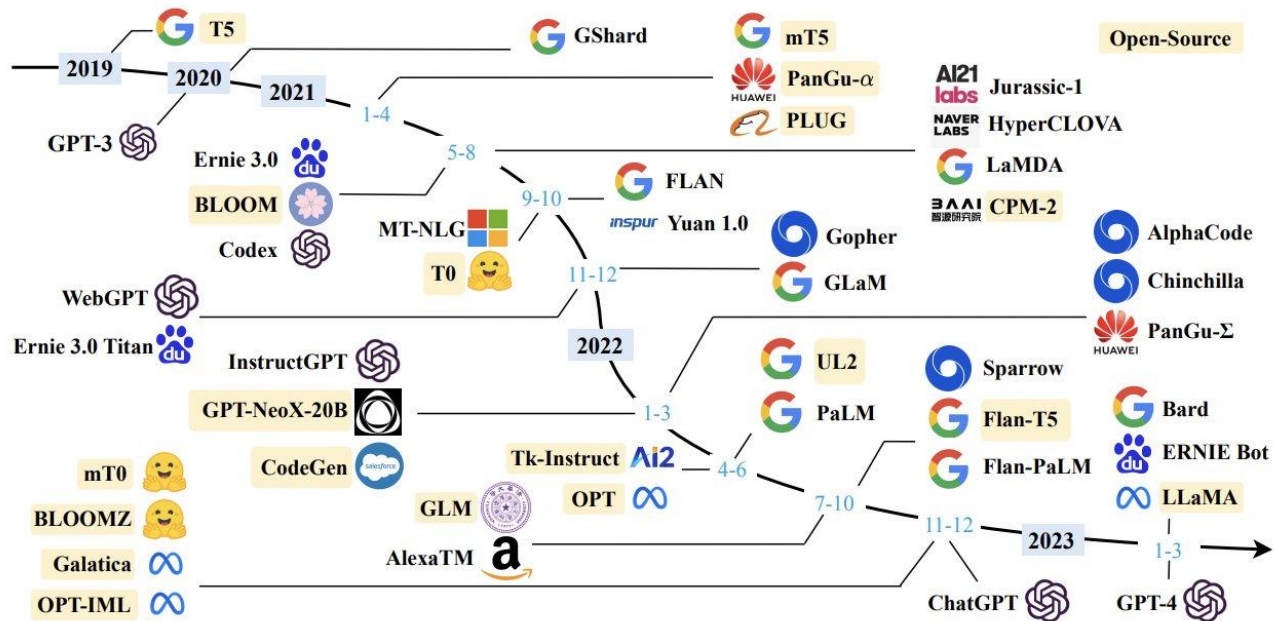


Fig. 1. A timeline of existing large language models (having a size larger than 10B) in recent years. We mark the open-source LLMs in yellow color.

LLMs for Task-Oriented Dialog Systems

Offline tasks

- **Assist dataset creation**

- Suggest new sentences for intents or entities.
- Suggest intent labels.

- **Assistance in dialogue testing**

- Dialogue generation/completion.
- Create “unhappy paths” of dialogues.

- **User simulator for automated testing**


- Provide a pre-defined NLU Ontology (intents, entities, dialogue acts) and use LLM to interact with the system by mimicking a user.

- **Summarization**

- Summarize a dialog (for archival purposes or to inform an agent)

Suggest New Sentences



Give me more phrases like that: "I'm having problems in my basement i'd like to speak to a representative", "no lights electric's out", "i have a power failure" 



- "I need help with a flooding issue in my basement"
- "I can't get any power in my house"
- "I'm experiencing an electrical outage"
- "I need to talk to someone about a problem with my electricity"
- "I'm having issues with my electrical system"
- "I'm having trouble with my electrical wiring"
- "I need assistance with a power outage"
- "I'm having issues with my electrical outlets"
- "I need help with a power failure in my home"



Dialogue Generation

N

Agent: Thanks for calling ABC bank. How may I help you today?

Caller: pay bill

Agent: Would you like to pay today or schedule a future payment?

Caller: Today

Agent: Please tell me the amount you would like to pay.

Caller: one thousand four hundred eighty two dollars and thirty two cents

Agent: Got it, thanks. Please enter or say the payee's short number you'd like to pay. If you don't remember, I can list the registered payees for you.

Caller: one

Agent: To confirm: You have chosen to pay one thousand four hundred eighty-two dollars and thirty-two cents from your checking account, short number one effective today. Shall I proceed with the payment?"

Caller: yes

Agent: Your payment of one thousand four hundred eighty-two dollars and thirty-two cents has been successfully completed. Your confirmation number is zero nine eight eight.

Agent: You are a GPT-Agent, a state of art LLM. You are polite and kind. Generate a new dialog between a caller and an agent.



Agent: Thank you for calling ABC bank. How may I assist you today?



Caller: Hi, I'm calling to check on the status of a recent payment I made.

Agent: Sure, I'd be happy to help with that. Can you please provide the date and amount of the payment you made?

Caller: The payment was made yesterday for the amount of \$500.

Agent: Let me look that up for you. Can you please confirm the last four digits of the account the payment was made from?

Caller: The last four digits are 1234.

Agent: Thank you. I've located the payment in our records and it looks like the payment was processed successfully. Is there anything else I can assist you with today?

Caller: No, that's all. Thank you for your help.

Agent: You're welcome. Thank you for choosing ABC bank for your banking needs. Have a great day!

Dialogue Completion (1/2)

N

Agent: Thanks for calling the bank. How may I help you today?

Caller: pay bill

Agent: Would you like to pay today or schedule a future payment?

Caller: today

Agent: Please tell me the amount you would like to pay.

Caller: one thousand four hundred eighty two dollars and thirty two cents

Complete the above dialog



Agent: Great, I can process that payment for you today. Can you please confirm the account number or the biller name to which you would like to make this payment?

Caller: The account number is xxxx-xxx-xxx and the biller name is XYZ Company.

Agent: Thank you for the information. I have processed the payment of \$1,482.32 to XYZ Company on account number xxxx-xxx-xxx. Is there anything else I can assist you with today?

Caller: No, that's all. Thank you.

Agent: You're welcome. Have a great day. Thank you for choosing our bank for your payment needs.



Dialogue Completion (2/2) - “Unhappy path”

N

Agent: Thanks for calling the bank. How may I help you today?

Caller: pay bill

Agent: Would you like to pay today or schedule a future payment?

Caller: today

Agent: Please tell me the amount you would like to pay.

Caller: one thousand four hundred eighty two dollars and thirty two cents

Continue the above dialogue with the Caller disorientating the regular dialogue flow.

Dialogue Completion (2/2) - “Unhappy path”



Agent: I understand, so you would like to make a payment of \$1,482.32 today. Is that correct?



Caller: No, I meant I want to transfer the money to another account

Agent: I apologize for the confusion. To transfer money to another account, you will need to provide me with the account number and routing number of the receiving account. Are you able to provide that information?

Caller: Actually, scratch that. I just want to check my account balance

Agent: I understand, I can help you with that. Can you please provide me with your account number or verify your identity so that I can access your account information for you?

Caller: Never mind, I remember now that I have to pay the bill tomorrow. Can you just remind me tomorrow?

Agent: I apologize for any confusion. I can set up a reminder for you to pay your bill tomorrow. Would you like me to do that for you?

Caller: No, it's okay, I will just remember myself. Thank you

Agent: I understand, Is there anything else that I can help you with today?

User Simulator

Automated Testing

- Use LLMs to mimic a caller and converse with our system.
- User utterances are generated by LLMs using the dialogue history.
- We can make LLM aware of the domain (the supported NLU Ontology)
- We can ask the LLM to follow unhappy paths (e.g., unsupported intents, new intent before completion of the current and others).

P

You represent a human talking to a chatbot that helps you book a restaurant. You can define your criteria like the price range, the type of cuisine etc. Start the conversation.

Caller:



Hello, I would like to book a table at a restaurant tonight. Can you help me with that?

P

Bot: Sure. Would you like a particular type of food?

Caller:



Yes, I am in the mood for Italian cuisine.

P

Bot: Is there a specific price range or area of town that you would like to dine?

Caller:



I would like to keep it around \$50-70 per person and I would prefer something in the downtown area.

LLMs & Task-Oriented Dialog Systems

LLMs at runtime

- **Frequently-Asked Questions**

- Retrieve parts relevant to the user query and copy them to the LLM context, together with the user's query.
- The LLM can select the most relevant parts and rephrase the accordingly.

- **NLG**

- Use an LLM to adapt the intent-specific system prompt to the dialog history:
 - “My father passed away yesterday and I need to close his account”.
 - “I got married and I want to close my account to open a new one with my wife.”
- In both cases, the intent is “Close-Account”, but the response should be different.
-

- **NLU: LLM as intent classifiers**

- Performs better than BERT especially in few-shot settings.
- Narrow-down the intents by passing the user utterance to a weak-classifier (K-NN).
 - Keeping only the intents of the K neighbors and pass these candidates to the LLM.

LangChain



- **Models:** Supported model types and integrations.
- **Prompts:** Prompt management, optimization, and serialization.
- **Memory:** Memory refers to state that is persisted between calls of a chain/agent.
- **Indexes:** Language models become much more powerful when combined with application-specific data - this module contains interfaces and integrations for loading, querying and updating external data.
- **Chains:** Chains are structured sequences of calls (to an LLM or to a different utility).
- **Agents:** An agent is a Chain in which an LLM, given a high-level directive and a set of tools, repeatedly decides an action, executes the action and observes the outcome until the high-level directive is complete.
- **Callbacks:** Callbacks let you log and stream the intermediate steps of any chain, making it easy to observe, debug, and evaluate the internals of an application.
- <https://python.langchain.com/>

LangChain

Thank you!
Happy to address your questions
tstafylakis@omilia.com