

# SOTA in Text To Speech



Евгений Шуранов  
Huawei Principal Engineer  
ITMO associate professor  
PhD

[evgeniy.shuranov@huawei.com](mailto:evgeniy.shuranov@huawei.com)



ITMO UNIVERSITY



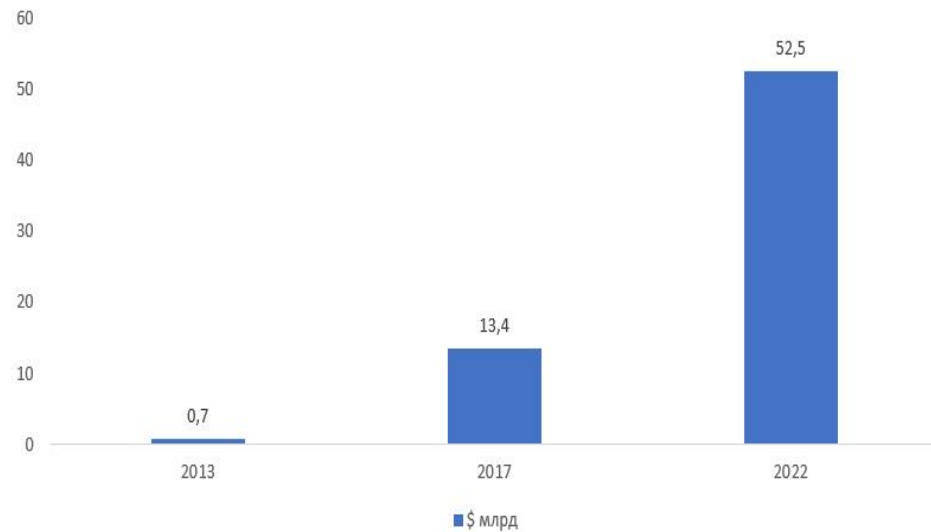
# Agenda

---

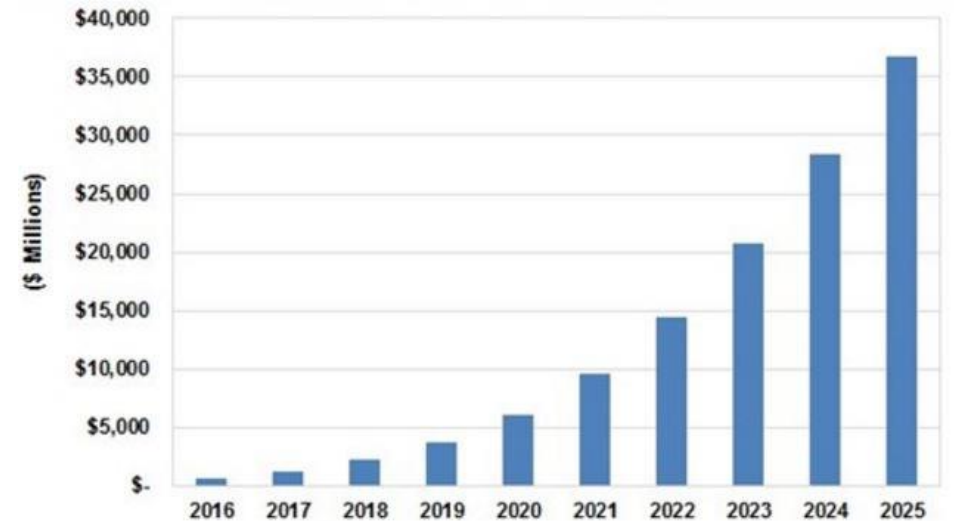
- Overview
  - Directions in speech processing
  - First steps
- SOTA TTS
- Example of interesting research in TTS
  - Multi-lingvo Multi-speaker Text-to-Speech
  - TTS improves speech recognition
  - Emotional Text-to-Speech system
- About our team and contacts

# Рост рынка ИИ

Объем мирового рынка технологий ИИ в период 2013-2022 гг.

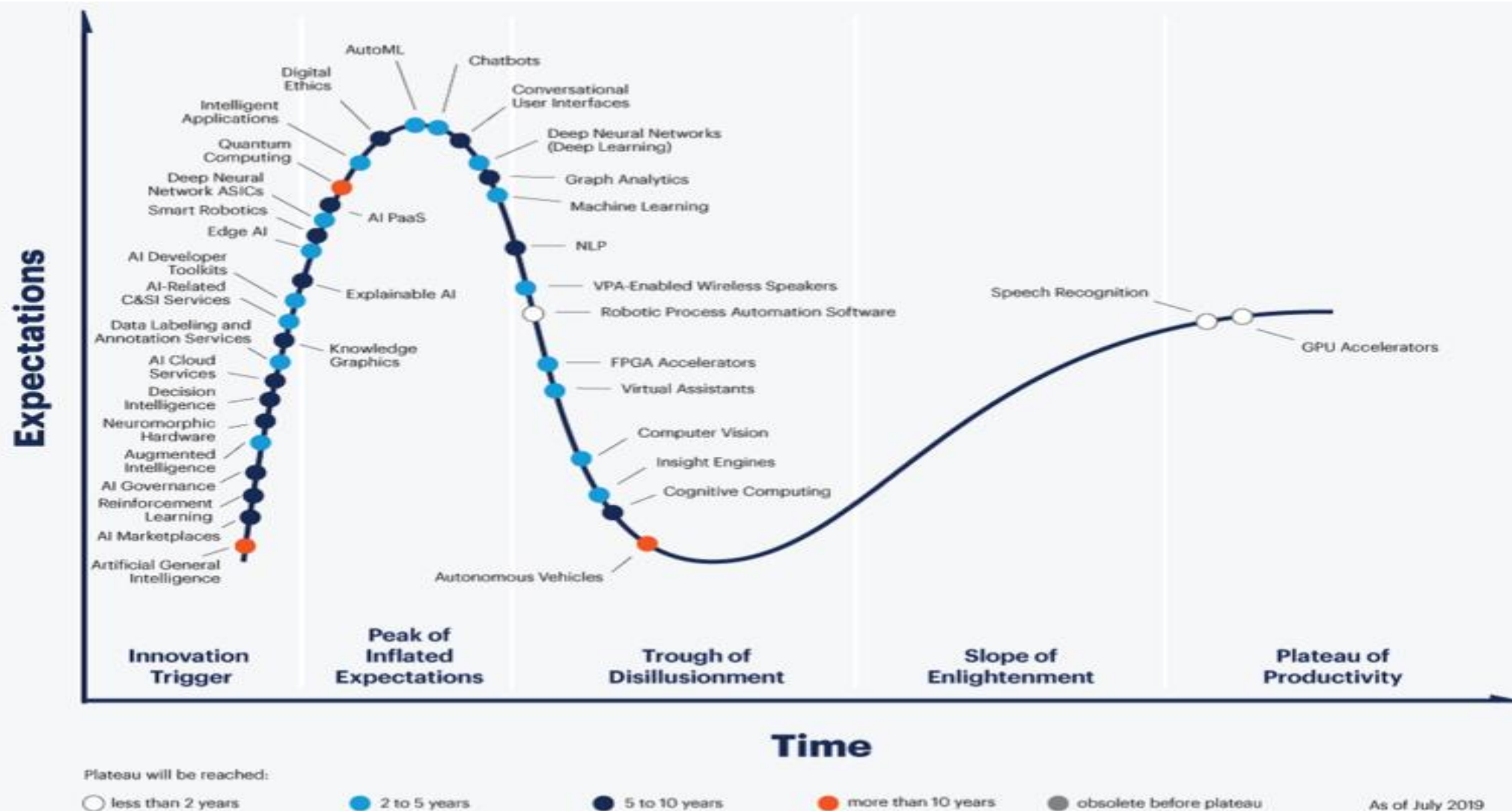


Доходы рынка ИИ 2016-2025. Данные Tractica



Рынок ИИ будет расти на 31% ежегодно -Frost & Sullivan

# Expectations



# How to start:

---

What you need	How to find
Papers	Arxiv / conferences / <b>challenges</b> /habr /telegram /slark and e.t.c.
Source examples	-----
Datasets	-----
GPU	google colab / universities

# Main events in speech technology

Name	Main domain	Type
Blizzard	TTS	Challenge
Voice Conversion Challenge	TTS / Voice Conversion	Challenge
Chime	ASR	Challenge
NIST	Voice recognition	Challenge
Interspeech	ASR++	Conference
Speecom	ASR++	Conference
Icasp	ASR++	Conference
EmotiW	Emo recognition	Challenge
Dcase	Event detection	Challenge
OMG	Emo recognition	Challenge

# First step errors

1 Math is always not enough!



KJ Cheetham ❄️ #FBPE 💎  
@kj\_cheetham

Читать

A maths meme that is actually funny rather than stupid:  
Solve carefully!  
 $230 - 220 \times 0.5 =$

You probably won't believe it but the answer is 5!

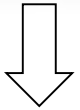
#maths

01:59 - 13 июл. 2019 г.

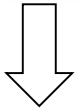
2 If you good at math remember about physics! (or business value ;) )

# What is TTS? Applications?

Any text to produce!



Text-to-Speech  
method

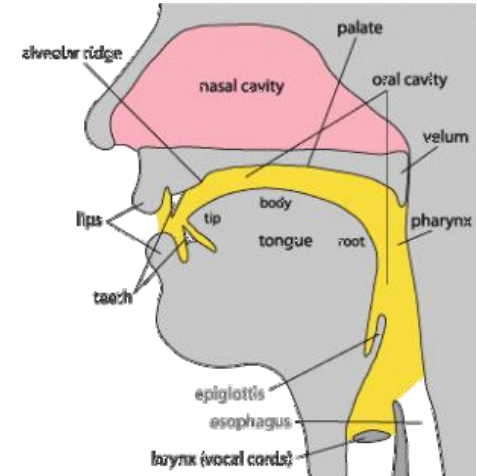


Intelligible speech



In 1779 the German scientist Christian Gottlieb Kratzenstein won the first prize in a competition announced by the Russian Imperial Academy of Sciences

and Arts for models he built of the human vocal tract that could produce the five long vowel sounds. In International Phonetic Alphabet notation: [a:], [e:], [i:], [o:] and [u:].

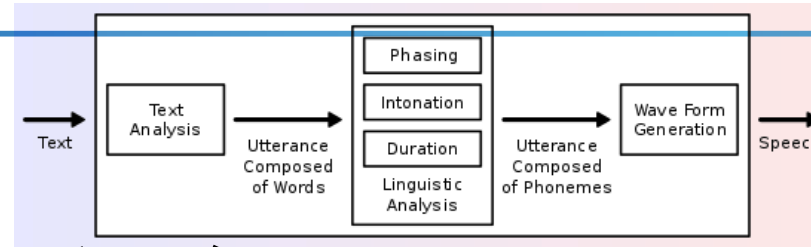




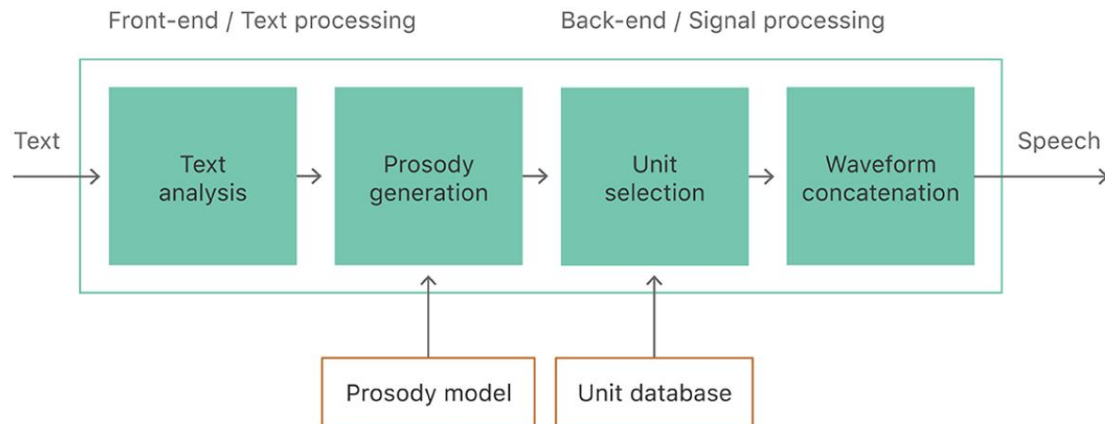
# Conventional TTS technologies

[7]

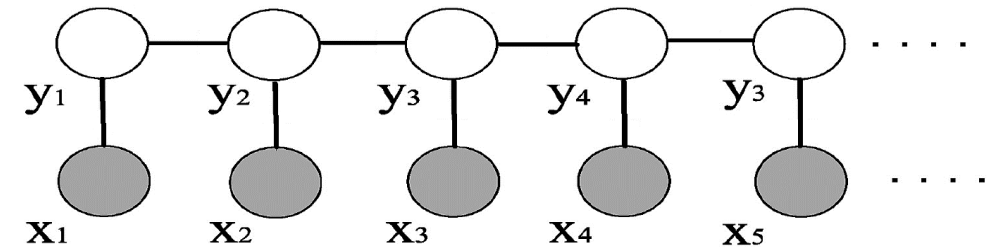
Overview of a typical TTS system



Concatenative unit selection



Statistical parametric Synthesis



Viterbi algorithm to produce speech waveform

# Conventional TTS technologies

[8]

## Front-End: Text preprocessing

### 1. Normalization

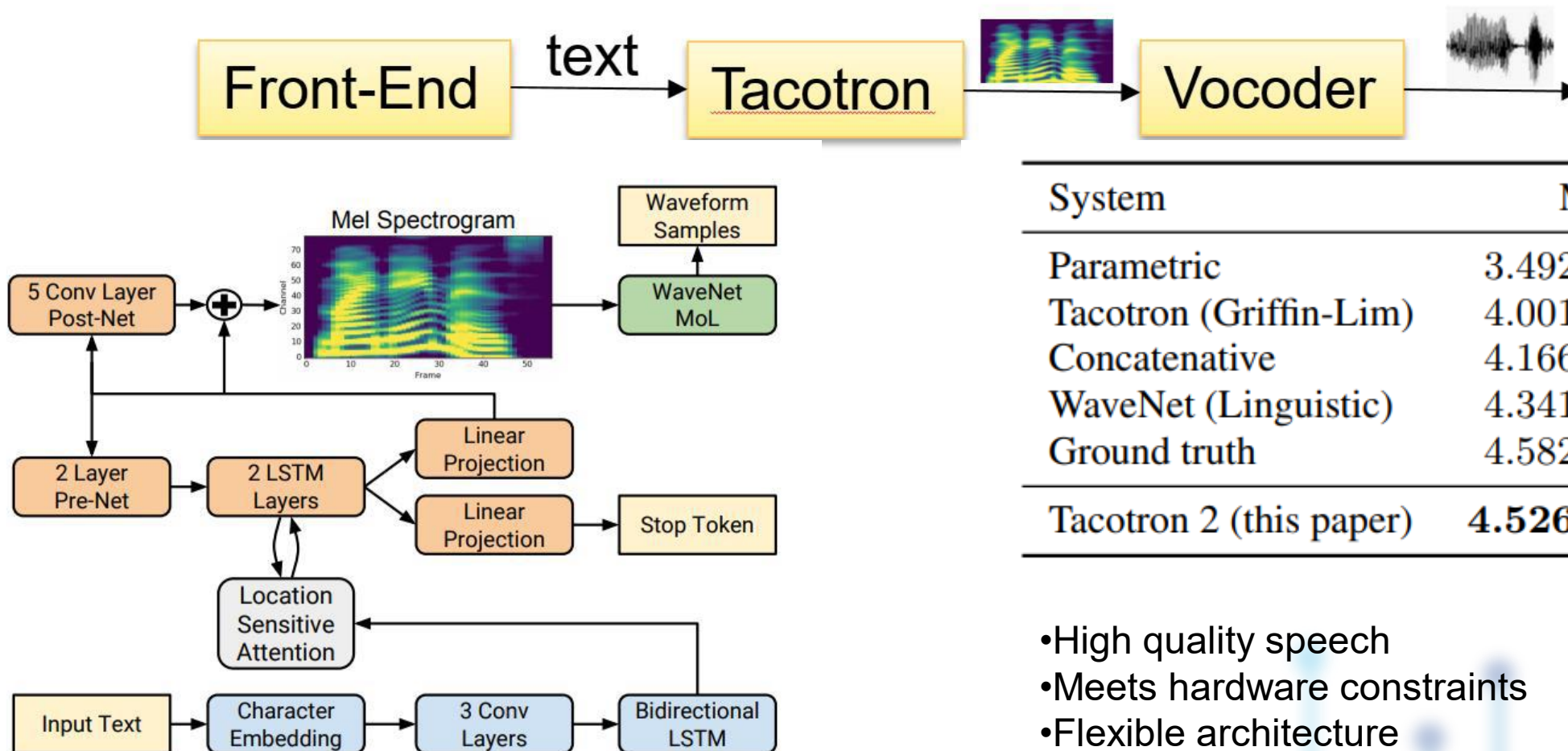
1. She was born on 1994
2. She was born on nineteen ninety four

### 2. Phonemization (IPA or Arpabet for English)

1. She was born on 1994 –
2. SH IY1 / W AA1 Z / B AO1 R N / AA1 N / W AH1 N / TH AW1 . Z AH0 N D / N AY1 N / HH AH1 N . D R AH0 D / AH0 N D / N AY1 N . T IY0 / F AO1 R
3. 你好, 你好吗? (nǐ hǎo nǐ hǎo ma)

alpha	EXPN	abbreviation	<i>adv, N.Y, mph, gov't</i>
	LSEQ	letter sequence	<i>CIA, D.C, CDs</i>
	ASWD	read as word	<i>CAT, proper names</i>
	MSPL	misspelling	<i>geogaphy</i>
	NUM	number (cardinal)	<i>12, 45, 1/2, 0-6</i>
	NORD	number (ordinal)	<i>May 7, 3rd, Bill Gates III</i>
	NTEL	telephone (or part of)	<i>212 555-4523</i>
	NDIG	number as digits	<i>Room 101</i>
N	NIDE	identifier	<i>747, 386, I5, pc110, 3A</i>
U	NADDR	number as street address	<i>5000 Pennsylvania, 4523 Forbes</i>
M	NZIP	zip code or PO Box	<i>91020</i>
B	NTIME	a (compound) time	<i>3-20, 11:45</i>
E	NDATE	a (compound) date	<i>2/2/99, 14/03/87 (or US) 03/14/87</i>
R	NYER	year(s)	<i>1998, 80s, 1900s, 2003</i>
S	MONEY	money (US or other)	<i>\$3-45, HK\$300, Y20,000, \$200K</i>
	BMONEY	money tr/m/billions	<i>\$3-45 billion</i>
	PRCT	percentage	<i>75%, 3-4%</i>
	SPLT	mixed or "split"	<i>WS99, x220, 2-car</i> (see also SLNT and PUNC examples)
M	SLNT	not spoken, word boundary	word boundary or emphasis character: <i>M.bath, KENT*RLTY, really_</i>
	PUNC	not spoken, phrase boundary	non-standard punctuation: "****" in <i>\$99,9K***Whites, "... in DECIDE... Year</i>
I			
S			
C	FNSP	funny spelling	<i>sllooooooww, sh*t</i>
	URL	url, pathname or email	<i>http://apj.co.uk, /usr/local, phj@tpt.com</i>
	NONE	should be ignored	<i>ascii art, formatting junk</i>

# SOTA TTS



System	MOS
Parametric	$3.492 \pm 0.096$
Tacotron (Griffin-Lim)	$4.001 \pm 0.087$
Concatenative	$4.166 \pm 0.091$
WaveNet (Linguistic)	$4.341 \pm 0.051$
Ground truth	$4.582 \pm 0.053$
<b>Tacotron 2 (this paper)</b>	<b><math>4.526 \pm 0.066</math></b>

- High quality speech
- Meets hardware constraints
- Flexible architecture

# SOTA. Vocoder evolution.

Vocoder	Year	Type of network	Speed	Training time
Wavenet	2016	Fully-convolutional	Extremelyslow (Open versions take several minutes to produce 1 second)	Weeks
WaveRNN	2018	Fullyrecurrent	4xfaster than realtimeon GPU	Around 1 week (according to open implementations)
LPCNet	2018	Convolutional+Reccurent+ Algorithmic(LPC)	4x faster than realtimeon CPU	1-2 days
MelGAN	2019	Fully-convolutional	20x faster than realtimeon CPU	3-4days

# Frontend

---

- Languages are very different

Hello Sota Machine learning School:



English



German(on English word)



Russian (on English word)



Russian

# Disentanglement

---

Speech data contains a lot of information:

Speech data = Text+ Speaker+ Prosody+ Recording conditions

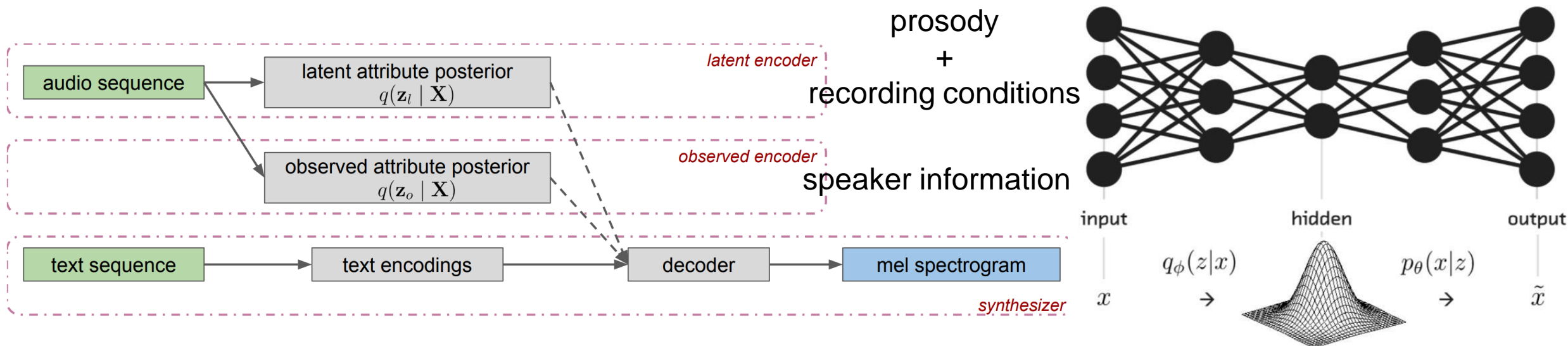
Some of them are present in the annotation and some are not.

- random or flat prosody,
- some level of noise during synthesis if recording artifacts were present in the data.

Aim:

1. We want to control all components of the speech during synthesis.
2. We want to train TTS without additional mark-up.

# Disentanglement



## GMVAE:

- uses hierarchical model for latent variables
- most of the component has interpretable role (noise/speed/pitch/accent etc) –easy control

## VAE:

- Latent space is not interpretable: Separate dimensions of the vector have no meaning
- Due to the above control (in practice) is problematic



# Multi-lingvo Multi-speaker Text-to-Speech

## Target:

Complexity: <50% bigger than single speaker model

Size: <50% bigger than single speaker model

Quality: same MOS as single speaker model

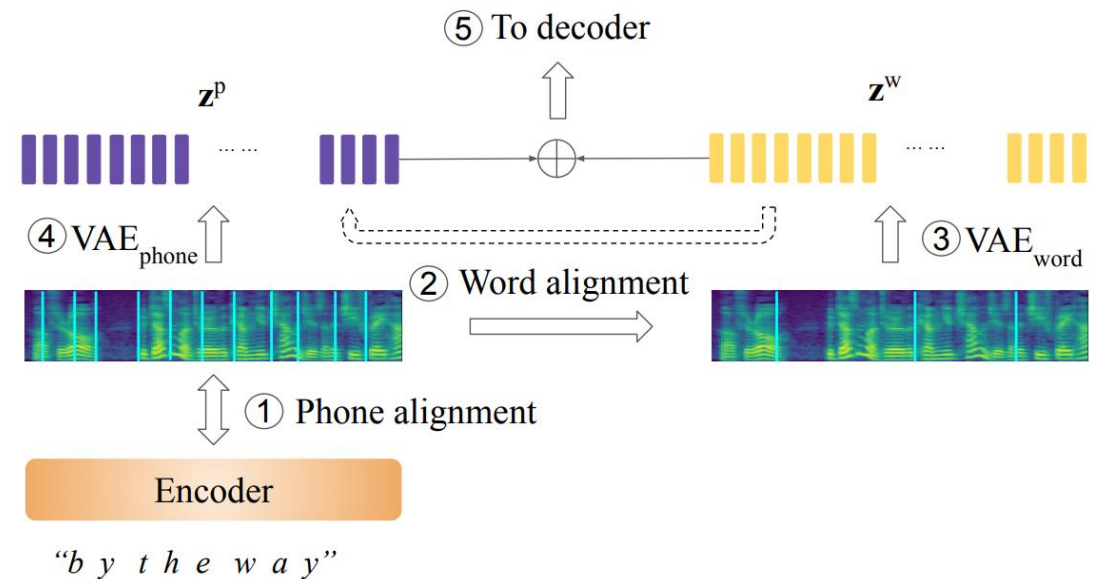
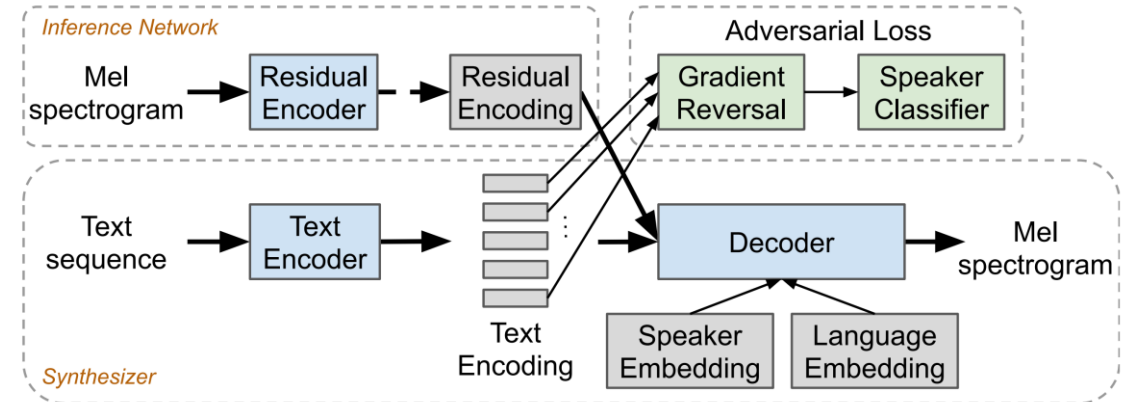
## Challenges:

Speaker/accent/language disentanglement

**NO product-ready open technology**

## Key technology:

VAE + RNN combination





# TTS improves speech recognition

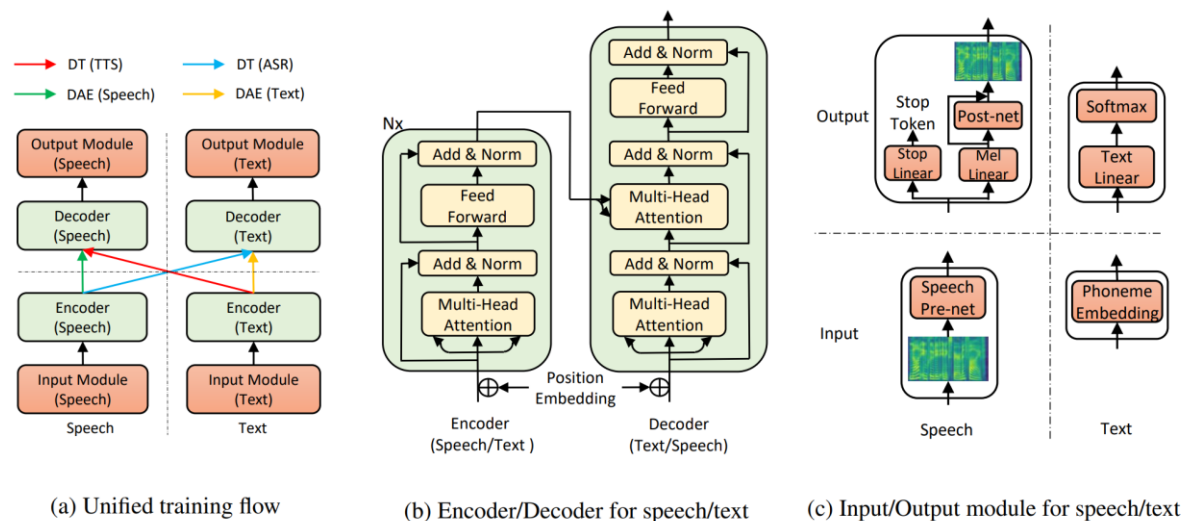


Figure 1. The overall model structure for TTS and ASR. Figure (a): The unified training flow of our method, which consists of a denoising auto-encoder (DAE) of speech and text, and dual transformation (DT) of TTS and ASR, both with bidirectional sequence modeling. Figure (b): The speech and text encoder and decoder based on Transformer. Figure (c): The input and output module for speech and text.

Method	MOS (TTS)	PER (ASR)
GT	4.54	-
GT (Griffin-Lim)	3.21	-
Supervised	3.04	2.5%
Pair 200	Null	72.3%
Our Method	2.68	11.7%

Table 1. The comparison between our method and other systems on the performance of TTS and ASR.

Paired Data	100	200	300	400	500
PER (ASR)	64.2%	11.7%	8.4%	5.2%	4.4%
MOS (TTS)	Null	2.45	2.49	2.64	2.78

Table 3. The PER on ASR with different amount of paired data for our method.

Authors used mutual training of unified transformers architecture for TTS, ASR.

Authors achieved 99.84% in terms of word level intelligible rate and 2.68 MOS for TTS, and **4.4% PER for ASR with just 500 paired data on LJSpeech** dataset. They stated to prove that it is possible to train ASR with close to SOTA performance only with 500 audio clips (~1 hour).

# Emotional Text-to-Speech system

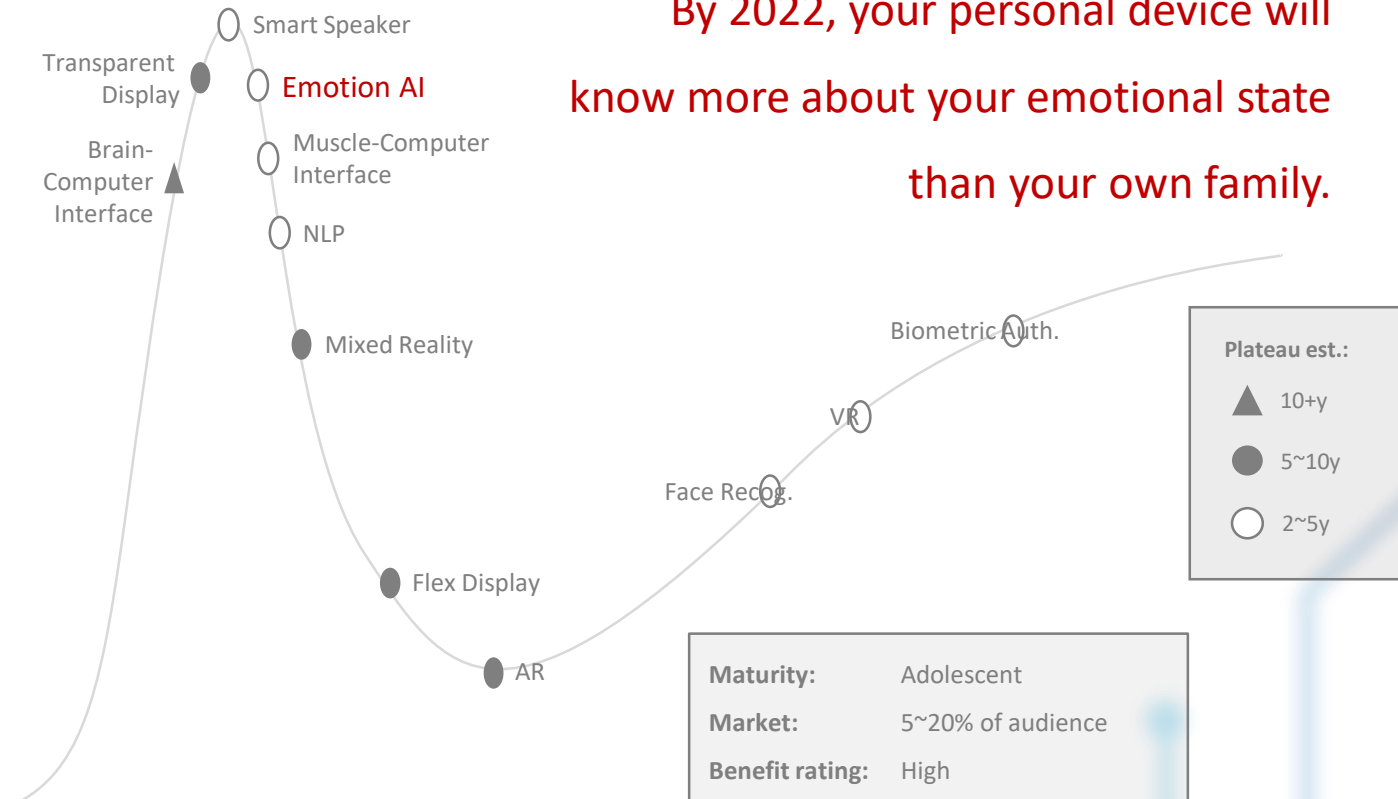
## 1) Speech Emotion Recognition (2020-2021)

- by voice
- by text

## 2) Emotional Text-to-Speech System (2020-2022)

### Challenges:

- poorly defined task
- not determined basic emotion list
- dataset transfer problems
- language transfer
- audio/text/video transfer
- actors/spontaneous speech transfer
- different labeling process
- model robustness and stability
- NO product-ready open technology



### Key technology:

pretrained features + LSTM/CNN combination

## Own research

### TTS improvement



### Voice cloning

Controlled TTS



Accent transfer



### ASR improvement



on device

### Event detection



Punctuation



## Cooperation research



Emotional TTS

Emotion recognition (Voice)

Emotion recognition (Text)

Voice-print technology



Anti spoofing



GANs approach (ASR + TTS)

# Goals of my team

---

- Find new perspective research directions in the speech domain
- Achieve result better than SOTA and find ways to implement this result for CBG products
- Make POC and production ready prototype
- Support delivery process of technology to production team

# Workshops & open days 2020 in SPb



01



СПбГУ



ПОЛИТЕХ

# Thank you for your attention

---



## Welcome to cooperation!



01

[https://docs.google.com/forms/d/e/1FAIpQLSfXzHrmdxdOKizMksoFi3IMSPoG\\_JMz4FH6QWYJ43IgUIdlQA/viewform](https://docs.google.com/forms/d/e/1FAIpQLSfXzHrmdxdOKizMksoFi3IMSPoG_JMz4FH6QWYJ43IgUIdlQA/viewform)



Contacts

Cooperation: [viktor.rakhmanov@huawei.com](mailto:viktor.rakhmanov@huawei.com)

HR: [nikulina.daria@huawei.com](mailto:nikulina.daria@huawei.com)