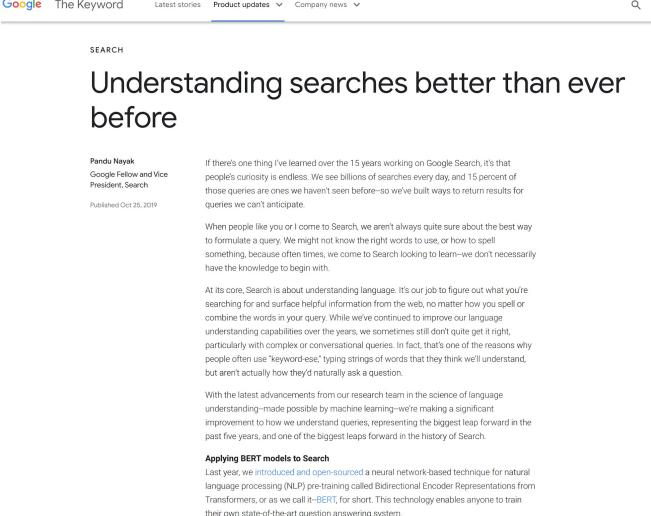


Achieving Low-latency Speech Synthesis at Scale

Sam Davis

Machine Learning Powers Products

A screenshot of a Google search results page for the query "The Keyword". The top result is from Google News, featuring a quote from Sundar Pichai: "Understanding searches better than ever before". Below it is a snippet from a blog post by Pandu Nayak, President of Search, dated October 25, 2019. The snippet discusses the evolution of search understanding over 15 years, mentioning BERT and its impact on query interpretation.

Sundar Pichai
Google CEO and President, Search

Published Oct 25, 2019

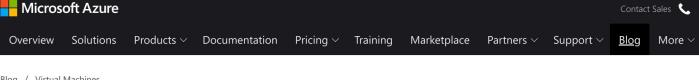
If there's one thing I've learned over the 15 years working on Google Search, it's that people's curiosity is endless. We see billions of searches every day, and 15 percent of those queries are ones we haven't seen before—so we've built ways to return results for queries we can't anticipate.

When people like you or I come to Search, we aren't always quite sure about the best way to formulate a query. We might not know the right words to use, or how to spell something, because often times, we come to Search looking to learn—we don't necessarily have the knowledge to begin with.

At its core, Search is about understanding language. It's our job to figure out what you're searching for and surface helpful information from the web, no matter how you spell or combine the words in your query. While we've continued to improve our language understanding capabilities over the years, we sometimes still don't quite get it right, particularly with complex or conversational queries. In fact, that's one of the reasons why people often use "keyword-ease," typing strings of words that they think we'll understand, but aren't actually how they'd naturally ask a question.

With the latest advancements from our research team in the science of language understanding—made possible by machine learning—we're making a significant improvement to how we understand queries, representing the biggest leap forward in the past five years, and one of the biggest leaps forward in the history of Search.

Applying BERT models to Search
Last year, we introduced and open-sourced a neural network-based technique for natural language processing (NLP) pre-training called Bidirectional Encoder Representations from Transformers, or as we call it—BERT, for short. This technology enables anyone to train their own state-of-the-art question answering system.

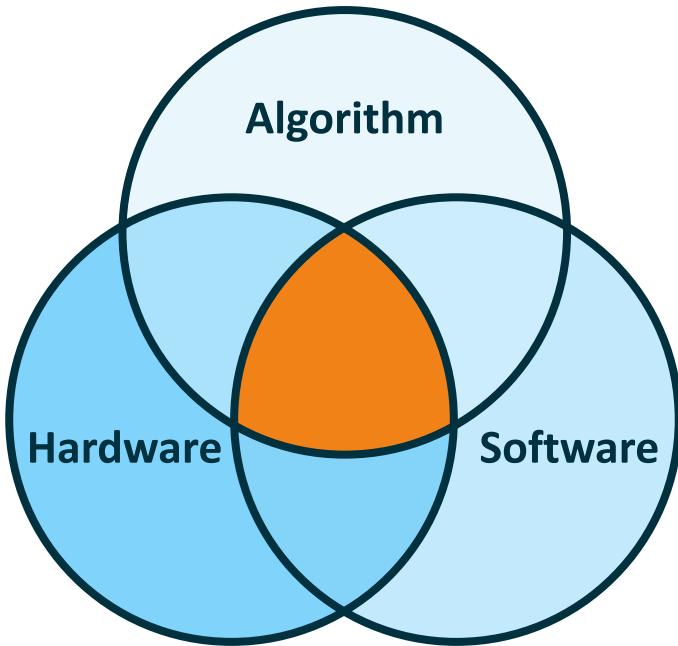
A screenshot of a Microsoft Azure search results page for the query "Bing delivers its largest improvement in search experience using Azure GPUs". The top result is from the Microsoft Azure blog, dated November 18, 2019, by Jeffrey Zhu, Program Manager, Bing Platform. The snippet highlights the use of Azure GPUs to improve search results by better understanding user intent and providing more relevant, contextual results.

Jeffrey Zhu, Program Manager, Bing Platform

Posted on November 18, 2019

Bing delivers its largest improvement in search experience using Azure GPUs

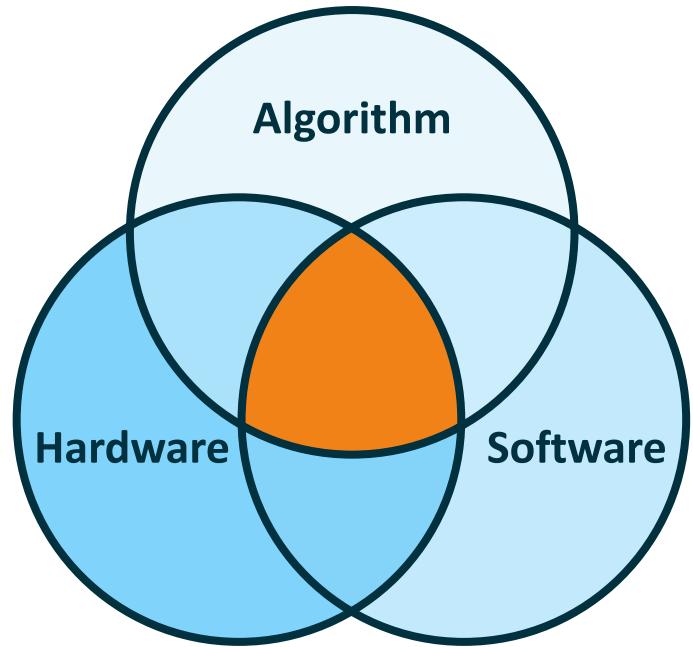
Powering the Future



Powering the Future

“The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of computation.”

- Rich Sutton, [The Bitter Lesson](#)



Talk Outline

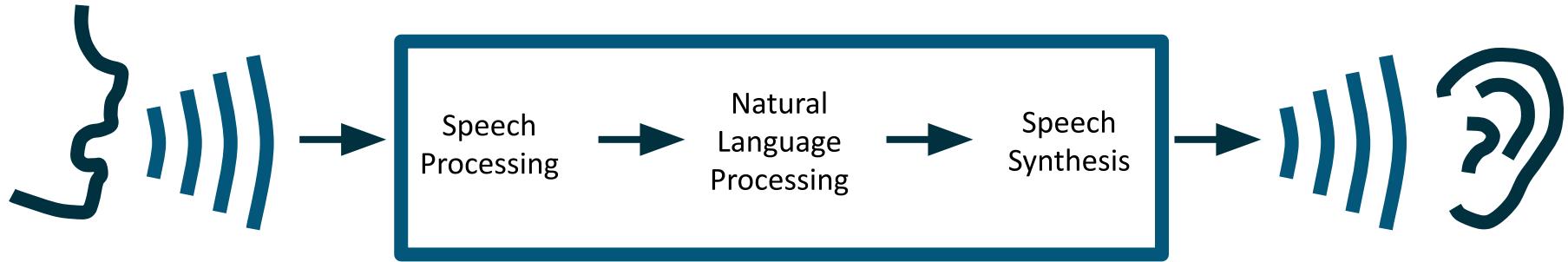
- Conversational AI Overview
- Speech Synthesis Solution Deep Dive
 - Algorithm
 - Hardware
 - Software
- Solution Demonstration and Performance



Myrtle.ai

Conversational AI

Conversational AI: What?



Conversational AI: Challenges

- **Users!**
 - Large volume of application requests
 - Each request uses multiple models
 - Each model has large throughput requirement
- **Quality**
 - Consistently return high quality results to satisfy users
 - Drives increase in model size and complexity
- **Latency**
 - Real-time interaction sets application latency bound to 100's ms
 - Model tail latencies critical to meet bound
 - Drives decrease in model size and complexity



Myrtle.ai

Solution

Speech Synthesis at Scale

Speech Synthesis Demo

Real Time Speech Synthesis on Intel® Stratix® 10 NX FPGA.

This demonstration showcases real time speech synthesis, running WaveNet on an Intel® Stratix® 10 NX FPGA. This deep neural network runs in BFP16 format on the AI Tensor Blocks and uses HBM memory.

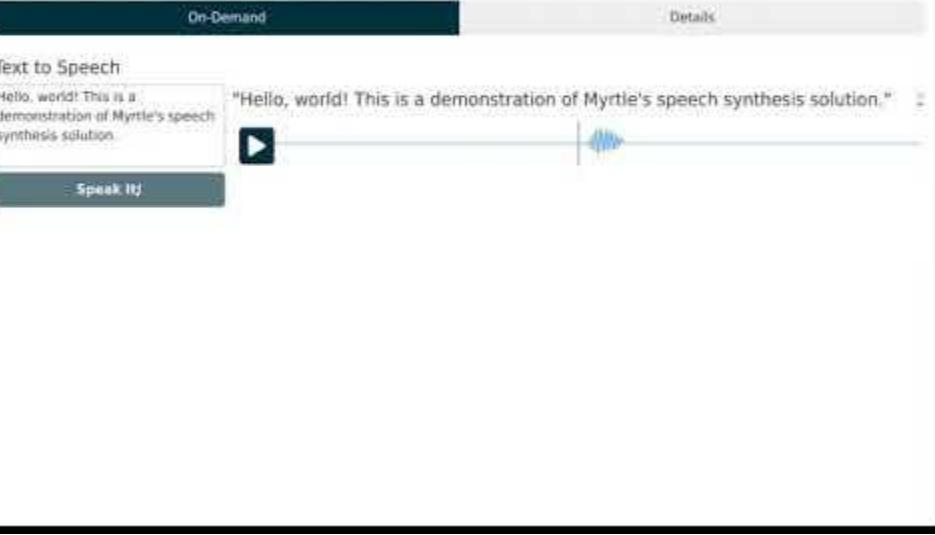
On-Demand Details

Text to Speech

Hello, world! This is a demonstration of Myrtle's speech synthesis solution.

Play

Speak It!





Myrtle.ai

Algorithm

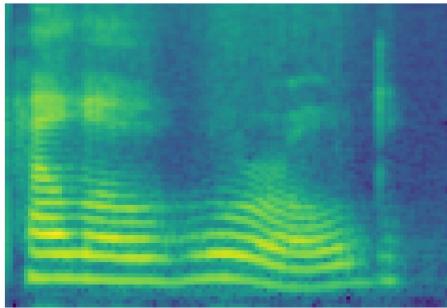
Neural Speech Synthesis

Neural Speech Synthesis Overview

Hello, world!



Tacotron 2

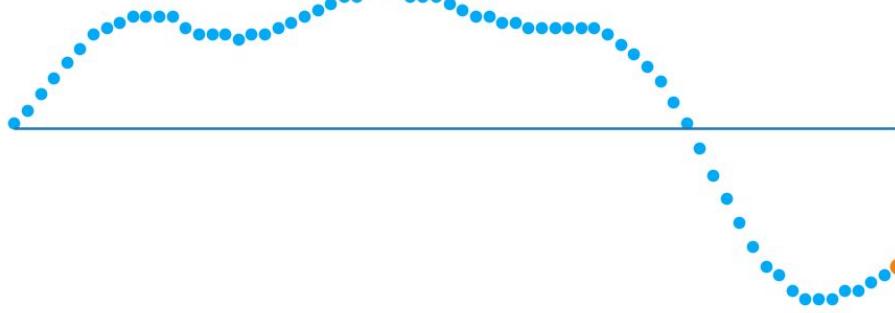


WaveNet

WaveNet: A Generative Model for Raw Audio

WaveNet models:

$$p(y_t | y_{<t}, f)$$



y_t Audio sample at step t. e.g. for 16-bit audio: y_t is in $[0, 2^{16})$

$y_{<t}$ All previous audio samples

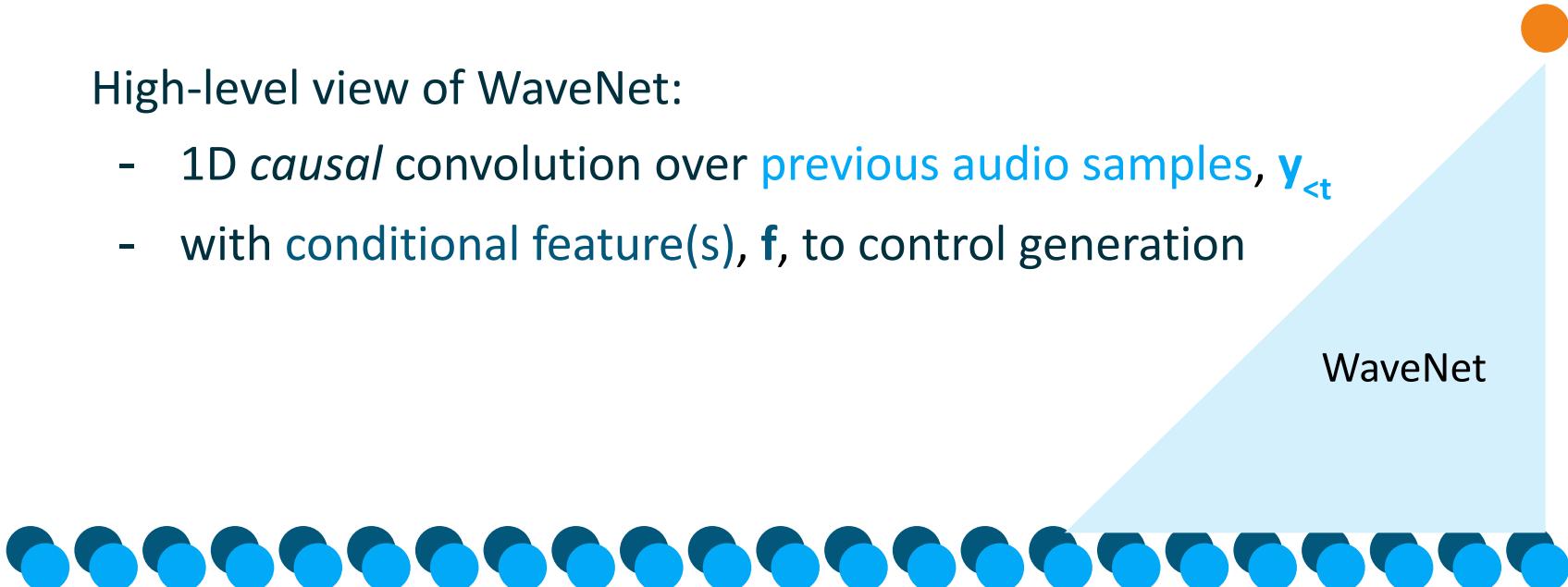
f Feature(s) to condition on. e.g. spectrogram

WaveNet: Model Overview

$$p(\mathbf{y}_t \mid \mathbf{y}_{<t}, \mathbf{f})$$

High-level view of WaveNet:

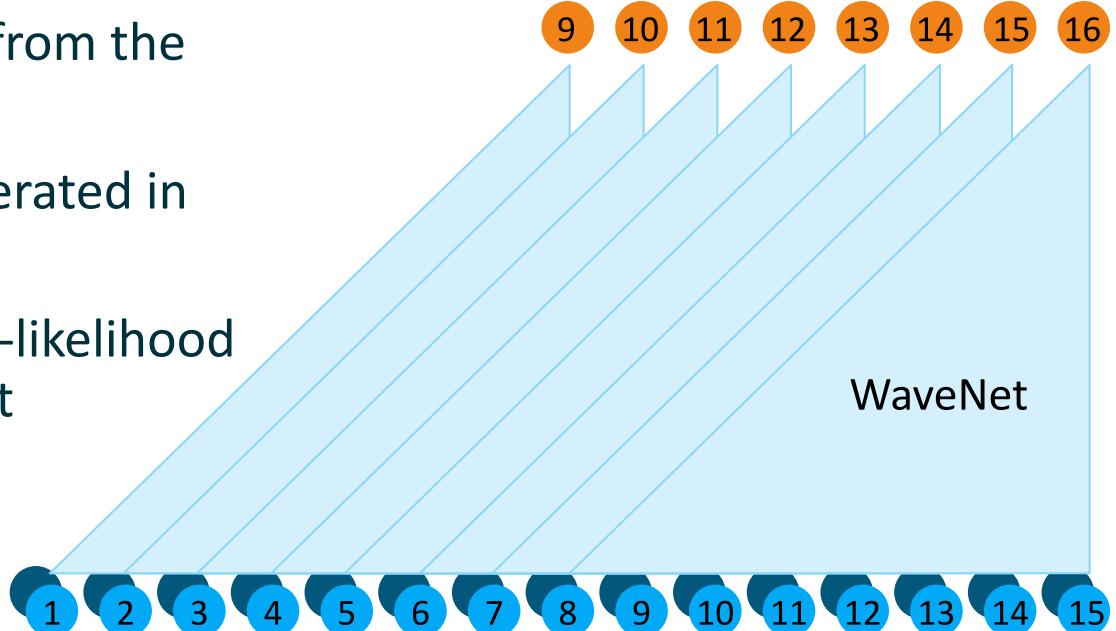
- 1D *causal* convolution over previous audio samples, $\mathbf{y}_{<t}$
- with conditional feature(s), \mathbf{f} , to control generation



WaveNet: Model Training

$$p(\mathbf{y}_t \mid \mathbf{y}_{<t}, \mathbf{f})$$

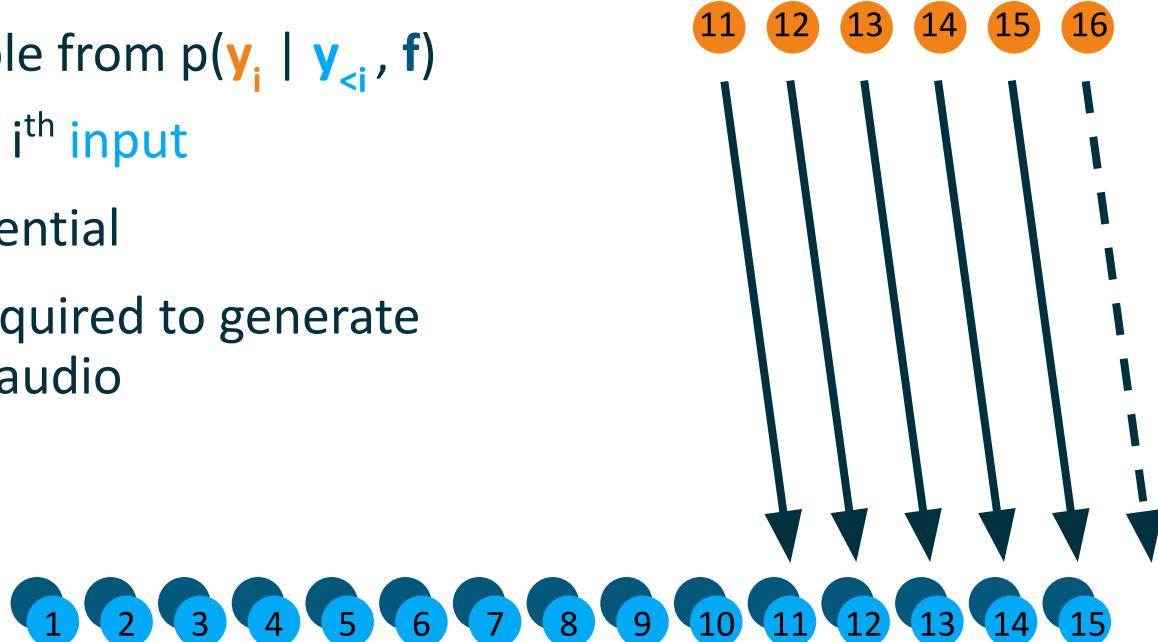
- Input and output are from the same audio clip
- All output can be generated in parallel
- Minimize negative log-likelihood using gradient descent



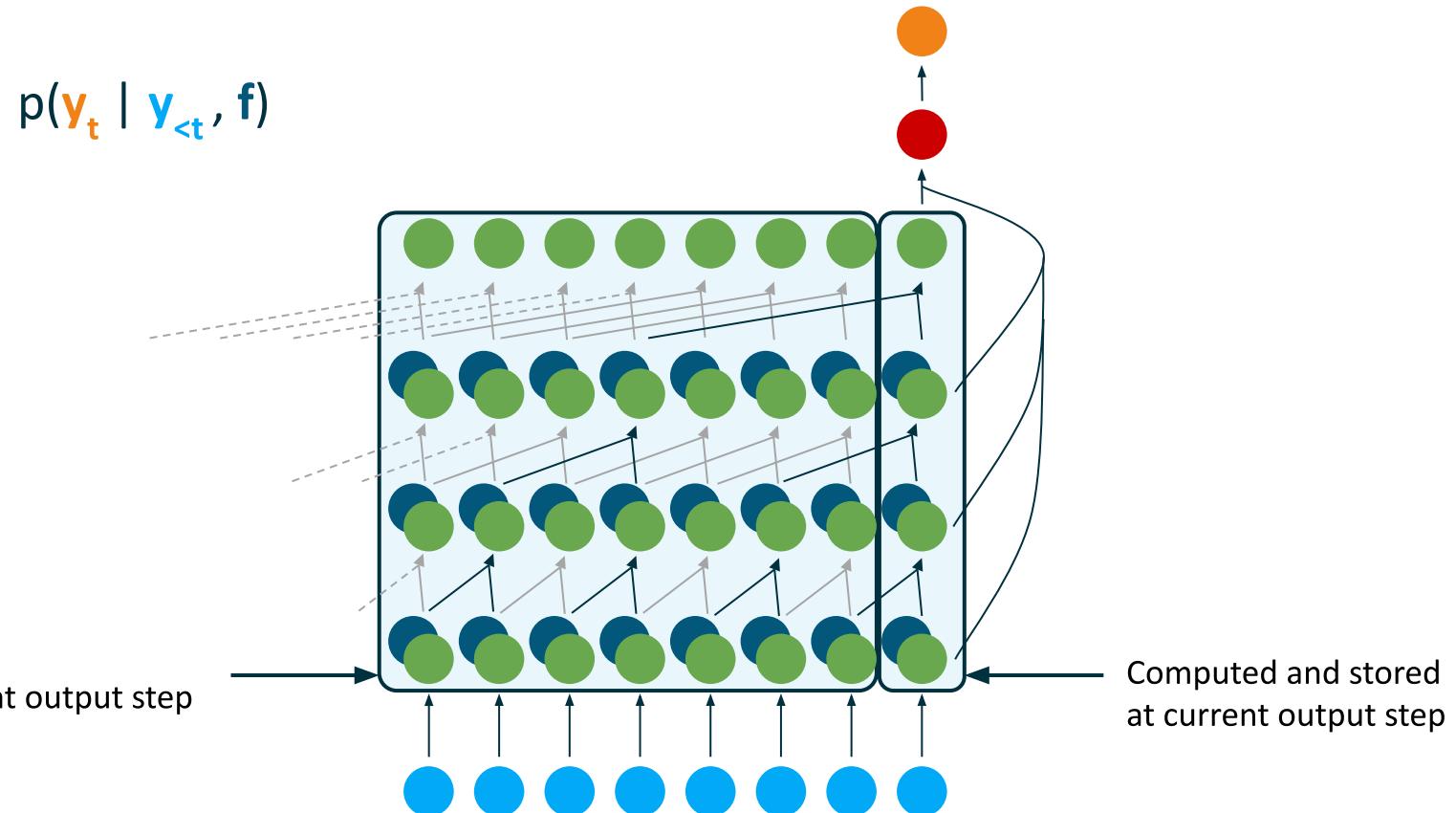
WaveNet: Model Inference

$$p(\mathbf{y}_t \mid \mathbf{y}_{<t}, \mathbf{f})$$

- i^{th} output is a sample from $p(\mathbf{y}_i \mid \mathbf{y}_{<i}, \mathbf{f})$
- i^{th} output becomes i^{th} input
- Generation is sequential
- 16000 iterations required to generate 1 second of 16kHz audio



WaveNet: Model Architecture



WaveNet: Model Architecture

$$p(y_t | y_{<t}, f)$$

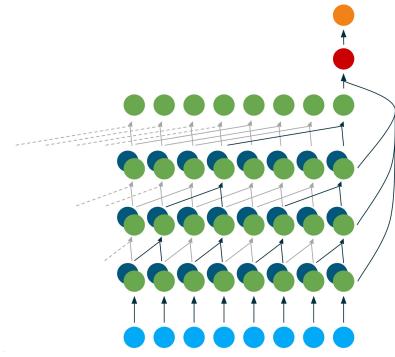
$$h_\theta = \text{Embed}(y_{<t})$$

$$z_k = \text{GatedActivationUnit}(h_{k-1}, f)$$

$$h_k = h_{k-1} + \text{Linear}(z_k)$$

$$o = \sum \text{Linear}(z_k)$$

$$p(y_t | y_{<t}, f) = (\text{Softmax} \circ \text{Linear} \circ \text{ReLU} \circ \text{Linear} \circ \text{ReLU})(o)$$



GatedActivationUnit(h_k, f) =
 $\tanh(\text{DilatedConv1D}(h_k) + f) \odot \sigma(\text{DilatedConv1D}(h_k) + f)$

WaveNet: Key Characteristics

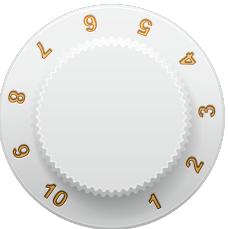
- Autoregressive: sequential generation
- Model backbone: deep stack of dilated convolutions

Quality Metric: Mean Opinion Score

How natural (i.e. human-sounding) is this recording?

RATING	LABEL	DESCRIPTION
1	Bad	Completely unnatural speech
2	Poor	Mostly unnatural speech
3	Fair	Equally natural and unnatural speech
4	Good	Mostly natural speech
5	Excellent	Completely natural speech

Quality Dials?



Model
Depth

Layer
Dimensions

Audio
Frequency

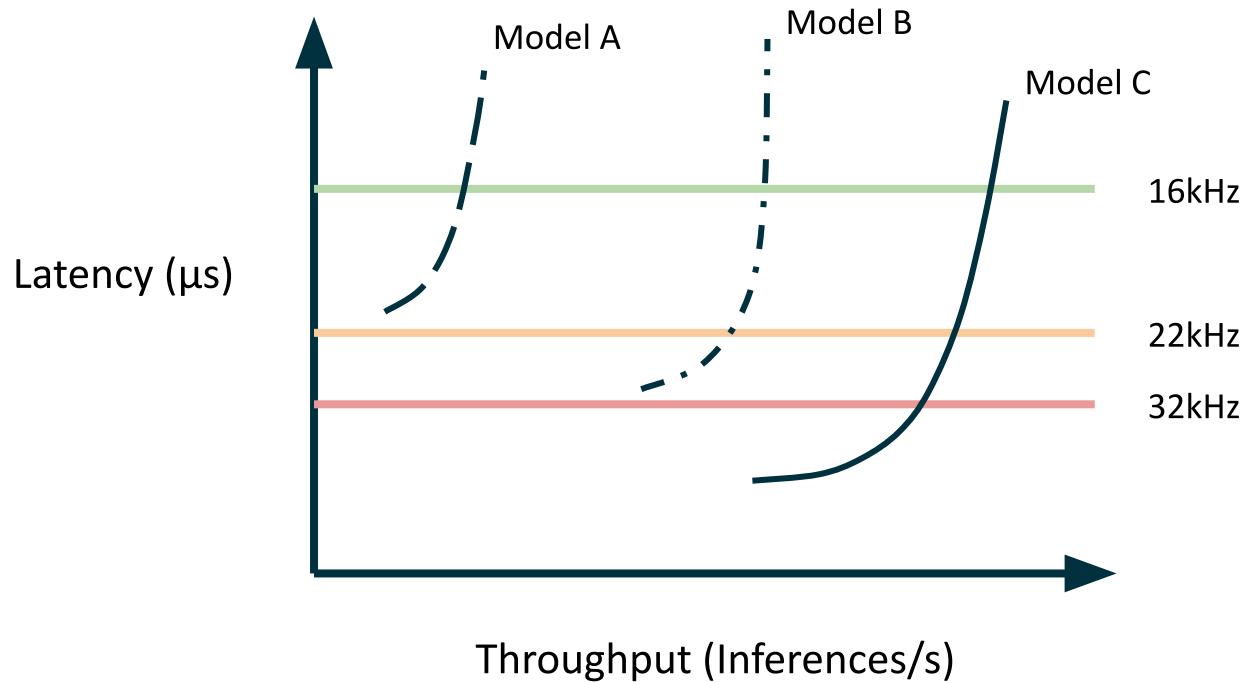


Myrtle.ai

Hardware

Intel® Stratix® 10 NX FPGA

Performance Metric: Latency-bound Throughput



Hard latency constraints for the single forward pass of a model in order to generate X kHz audio in real-time

WaveNet on Existing Accelerators?

```
void at::native::__... [purple] void cudnn::detail::implicit_con... [blue] [purple] void cudn... [blue] [blue] [purple] [purple] [purple] void cudnn::det... [purple] [blue] void cudnn::de... [blue] [blue] [blue]
```



```
void at::native::__... [purple] void cudnn::detail::implicit_con... [blue] [purple] void cudn... [blue] [blue] [purple] [purple] [purple] void cudnn::det... [blue] [blue] void cudnn::de... [blue] [blue] [blue]
```

$$z_k = \text{GatedActivationUnit}(h_{k-1}, f)$$
$$h_k = h_{k-1} + \text{Linear}(z_k)$$

Intel® Stratix® 10 NX FPGA

- **FPGA Architecture**
 - Adapt to rapidly changing model architectures and components
- **HBM2**
 - Meet neural network parameter and activation memory bandwidth requirements
- **AI Tensor Blocks**
 - Meet neural network compute requirements (OPS)

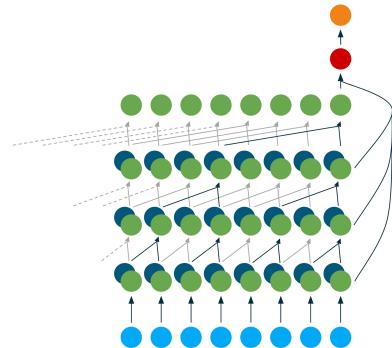


FPGA Architecture: Efficient Sampling

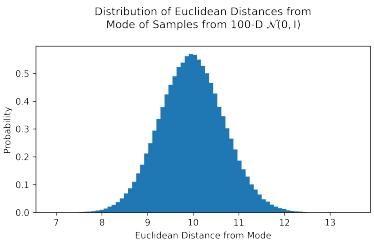
$\mathbf{l} = (\text{Linear} \circ \text{ReLU} \circ \text{Linear} \circ \text{ReLU})(\mathbf{o})$

$p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f}) = \text{Softmax}(\mathbf{l})$

$\mathbf{y}_t = \text{sample}(p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f}))$



- Softmax(\mathbf{l}) and sample($p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f})$) is expensive (logic/area, cycles)
- Alternative: Always “sample” mode of $p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f})$?
 - Cheap: select index of largest value in \mathbf{l}
 - Problem: Quality metric, MOS, is poor
 - Why? Typicality!



FPGA Architecture: Efficient Sampling

$\mathbf{l} = (\text{Linear} \circ \text{ReLU} \circ \text{Linear} \circ \text{ReLU})(\mathbf{o})$

$p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f}) = \text{Softmax}(\mathbf{l})$

$\mathbf{y}_t = \text{sample}(p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f}))$

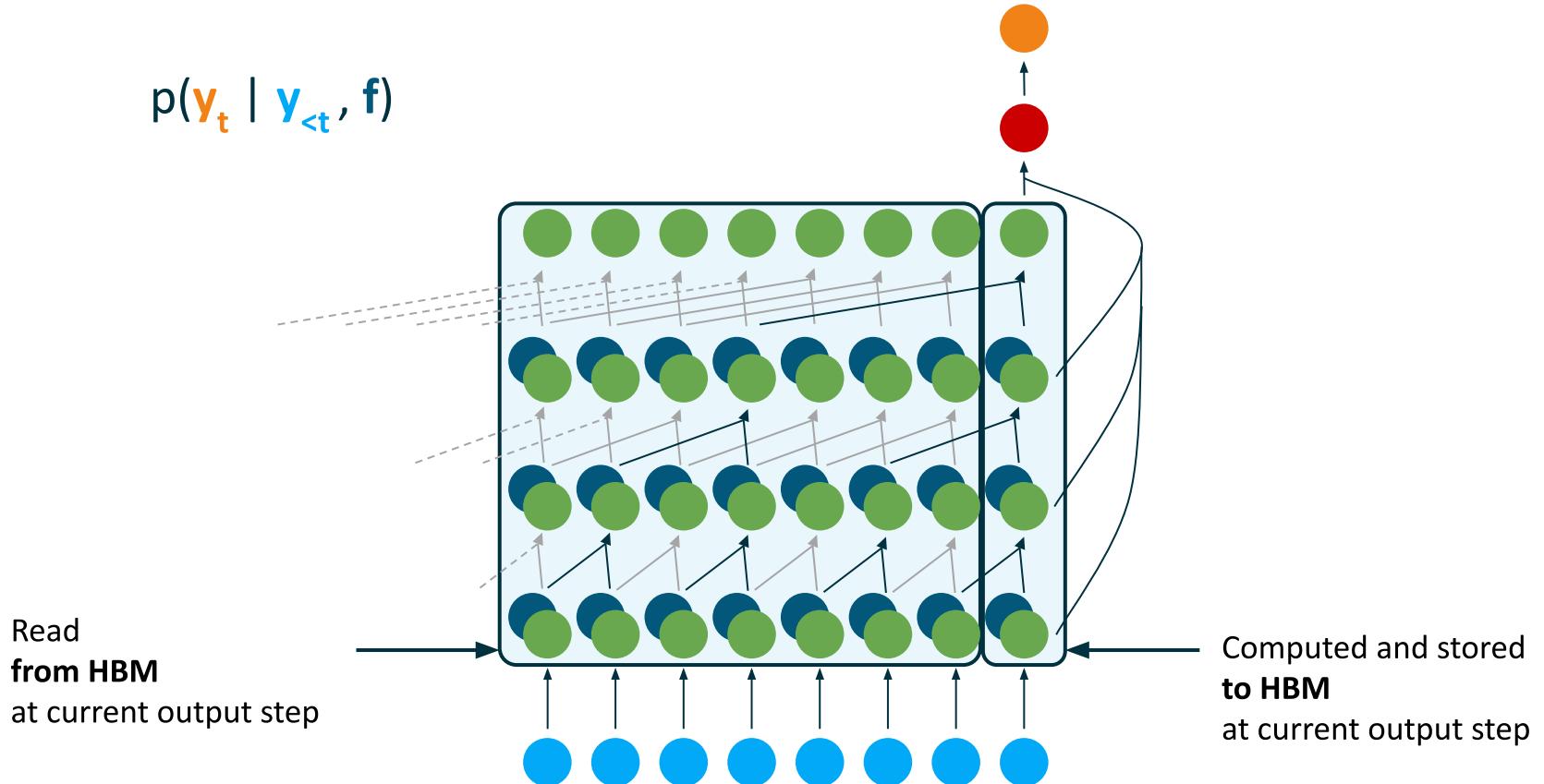
- Must sample, but how? Gumbel-max trick!

$\mathbf{y}_t = \text{sample}(p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f}))$
 $= \text{argmax}(\mathbf{l} + \mathbf{g})$

$\mathbf{g} \sim \text{Gumbel}(0, 1)$

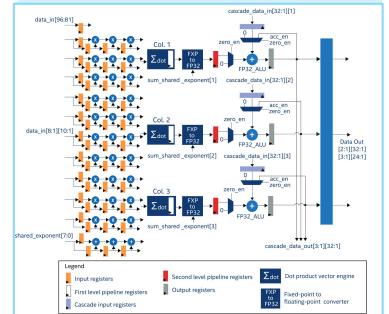
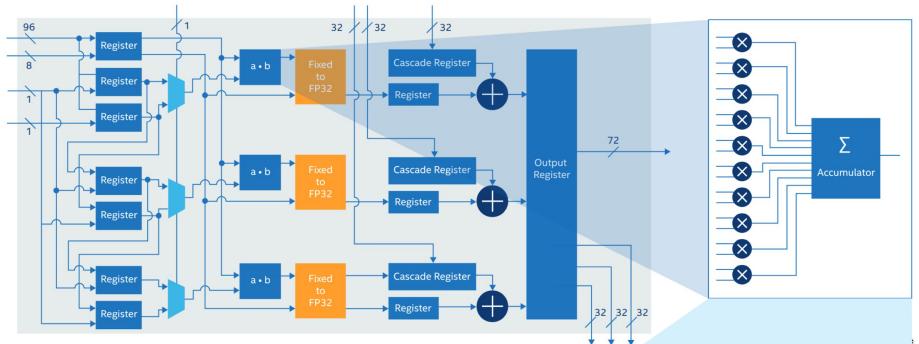
HBM2

$$p(\mathbf{y}_t \mid \mathbf{y}_{<t}, \mathbf{f})$$



AI Tensor Blocks

PRECISION	PERFORMANCE	EFFICIENCY
INT4	286 TOPS	2 TOPS/W
INT8	143 TOPS	1 TOPS/W
Block FP12	286 TFLOPS	2 TFLOPS/W
Block FP 16	143 TFLOPS	1 TFLOPS/W
@600 MHz Max Frequency		

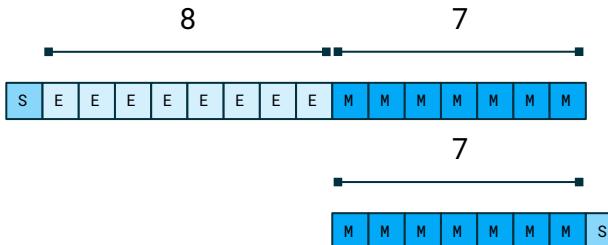


AI Tensor Blocks: Block Floating Point

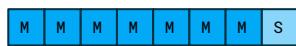
Single-precision Floating Point (FP32)



Brain Floating Point (bf16)



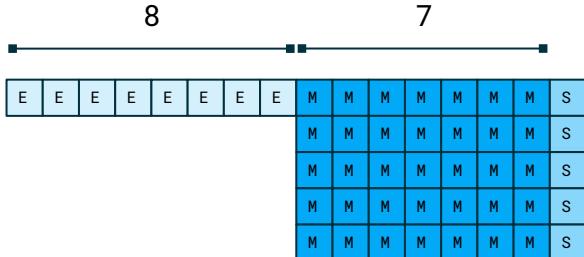
8-bit Integer (INT8)



Block Floating Point (BFP16)

or

“Microsoft Floating Point” (MSFP16)



Bounding block Size (e.g. 5)



Myrtle.ai



Software

Myrtle.ai MAU Accelerator

WaveNet: Performance?

Layer	Type	# Parameters		GOP/second audio	
		Per-layer	Total	Per-layer	Total
<i>Pre-processing Layers</i>					
Embedding	Embedding		30,720		-
<i>Repeated Layers</i>					
Dilation	Dilated Conv1d	57,840	925,440	1.84	29.49
Conditional	Conv1d	19,440	311,040	0.61	9.83
Residual	Conv1d	14,520	217,800	0.46	6.91
Skip	Conv1d	29,040	464,640	0.92	14.75
<i>Post-processing Layers</i>					
Out	Conv1d		61,440		1.96
End	Conv1d		65,536		2.09
<i>Total</i>			2,076,616		65.03

FPGA Accessibility

- Tooling and ecosystem limits usability and accessibility of FPGAs for ML and DL
- Deploying non-standard models requires deep understanding of full stack (ML -> FPGA/HW)
- Driving FPGA adoption for ML and DL requires improved platform, ecosystem, and tooling.

Myrtle.ai MAU Accelerator

- Myrtle.ai MAU Accelerator technologies and proven design techniques:
 - Deliver the right range of abstractions to all users across the stack
 - Enable Myrtle.ai to engage with partners to achieve 10x latency-bounded throughput on a wide range of workloads

Myrtle.ai MAU Accelerator: Architecture

- Configurable number of MAU Cores
- MAU Cores reconfigurable to make optimal use of target platform
- Abstraction reduces complexity of programming device and provides abstraction for ML engineers



Quantization - Quality Evaluation?

Last Layer Quality Impact

$\mathbf{z} = (\text{Linear2} \circ \text{ReLU} \circ \text{Linear1} \circ \text{ReLU})(\mathbf{o})$

$p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{f}) = \text{Softmax}(\mathbf{z})$

QAT	Linear2	Linear1	WaveNet MOS
No	FP32	FP32	3.930 ± 0.028
No	MSFP16	MSFP16	3.711 ± 0.029
Yes	MSFP16	MSFP16	3.823 ± 0.028
No	FP32	MSFP16	3.902 ± 0.027

Post-training Quantization

Category	Model	Float32	MSFP16
CNNs	Resnet-50	1.000 (75.26)	1.000
	Resnet-101	1.000 (76.21)	1.000
	Resnet-152	1.000 (76.58)	1.000
	Inception-v3	1.000 (77.98)	1.000
	Inception-v4	1.000 (80.18)	1.000
	MobileNet-V1	1.000 (70.90)	0.998
	VGG16	1.000 (70.93)	1.000
	VGG19	1.000 (71.02)	1.000
	EfficientNet-S	1.000 (77.61)	1.000
RNNs	EfficientNet-M	1.000 (78.98)	1.000
	EfficientNet-L	1.000 (80.47)	1.000
Transformers	Production-DR	1.000 (76.10)	1.000
	Production-DS	1.000 (73.10)	1.000
	BERT-MRPC	1.000 (88.39)	1.000
Transformers	BERT-SQuAD1.1	1.000 (88.45)	1.000
	BERT-SQuADv2	1.000 (77.23)	1.000

Part of Table 2 from: [Darvish Rouhani, Bita, et al. "Pushing the Limits of Narrow Precision Inferencing at Cloud Scale with Microsoft Floating Point." Advances in Neural Information Processing Systems 33 \(2020\)](#)



Myrtle.ai

Solution

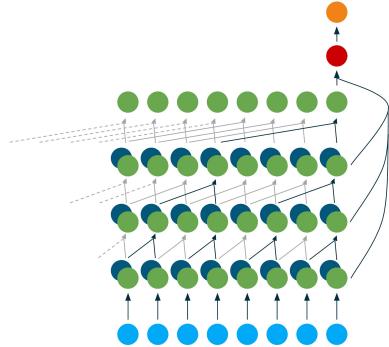
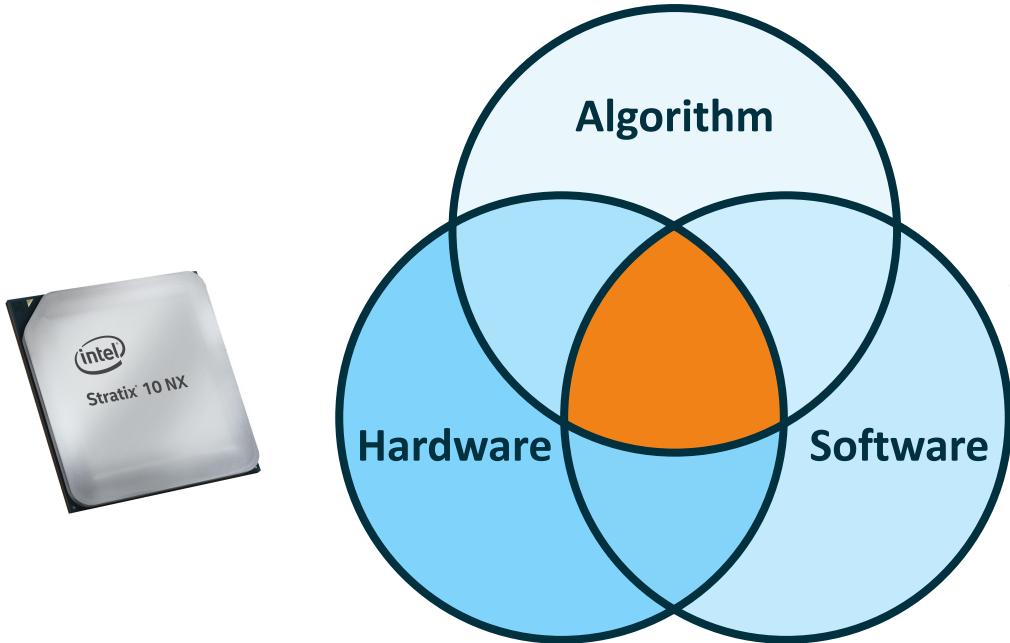
Speech Synthesis at Scale

using WaveNet

on Intel® Stratix® 10 NX FPGA

enabled by Myrtle.ai MAU Accelerator

Powering the Future



Solution Performance

Real Time Speech Synthesis on Intel® Stratix® 10 NX FPGA.



This demonstration showcases real time speech synthesis, running WaveNet on an Intel® Stratix® 10 NX FPGA. This deep neural network runs in BFP16 format on the AI Tensor Blocks and uses HBM memory. [Details](#)

On demand text to speech Text to speech stream
 2.29 TOPS 16 kHz 32.0 RTS 0 W 3.891 MOS (PTQ)

"and is very simple and legible, and unaffectedly designed for use: but it is by no means without beauty."



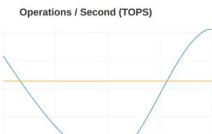
The characteristic Dutch type, as represented by the excellent printer Gerard Leew, is very pronounced and uncompromising Gothic.

some of which -- as, e.g., that of Jacobus Rubeus or Jacques le Rouge -- is scarcely distinguishable from his.

produced the block books, which were the immediate predecessors of the true printed book,

but his letters, though uninteresting and poor, are not nearly so gross and vulgar as those of either the Italian or the Frenchman.

that the forms of printed letters should follow more or less closely those of the written character, and they followed them very closely.



AUDIO FREQUENCY	CONCURRENT VOICE CHANNELS	
	MYRTLE.AI WAVENET	NV-WAVENET
16 kHz	256	32
24 kHz	160	16
32 kHz	128	8

More Information

- Solution content:
 - [Intel Blog](#)
 - [Whitepaper](#)
 - [Demo Video](#)
- Interactive solution [website](#)



Myrtle.ai

Thank You

www.myrtle.ai



[Sam Davis](#)



+44 1223 967248



sam@myrtle.ai