# Developing and Running Neural Audio in Constrained Environments

Carter Huffman
(CTO / Co-Founder, Modulate)
Brendan Kelly
(Research Engineer, Modulate)

# Objectives

1. Practical lessons and tools we have picked up to effectively train speech synthesis models.

2. Running these models real-time in gaming environments.

# Agenda

1. Overview of Neural Audio models
2. Challenges for Development
3. Challenges for Runtime
4. Epilogue

# Agenda

1. **Overview of Neural Audio models**
2. Challenges for Development
3. Challenges for Runtime
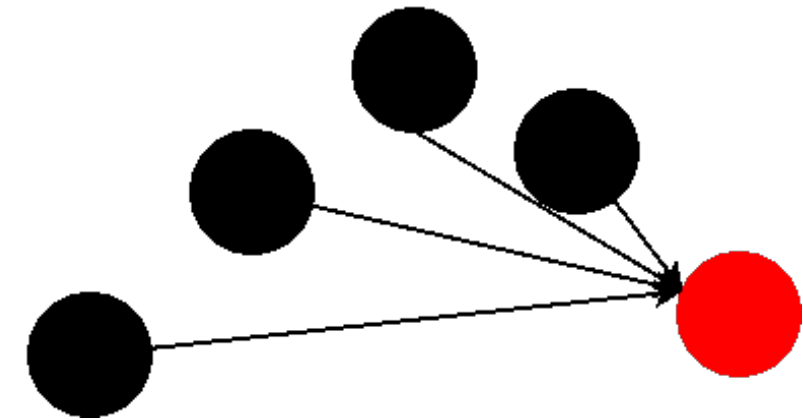4. Epilogue
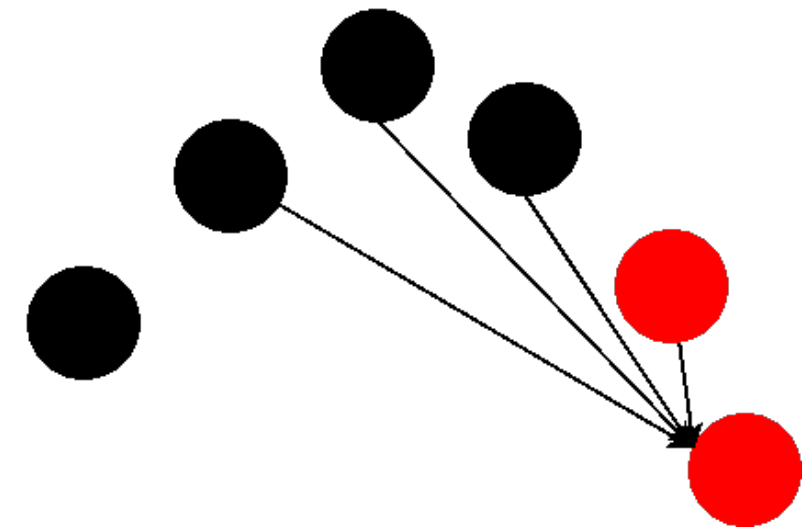
# Autoregressive Approach
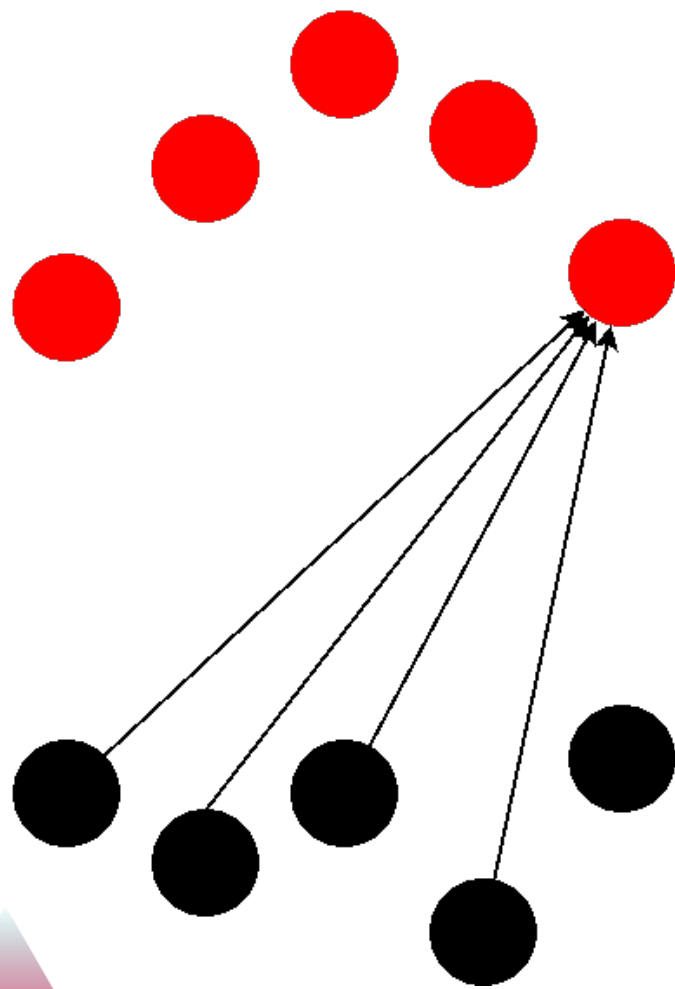
Parametric:

Autoregressive:

Predicting Sample n

Predicting Sample n+1

Reference: [a]

# Parallelized approach
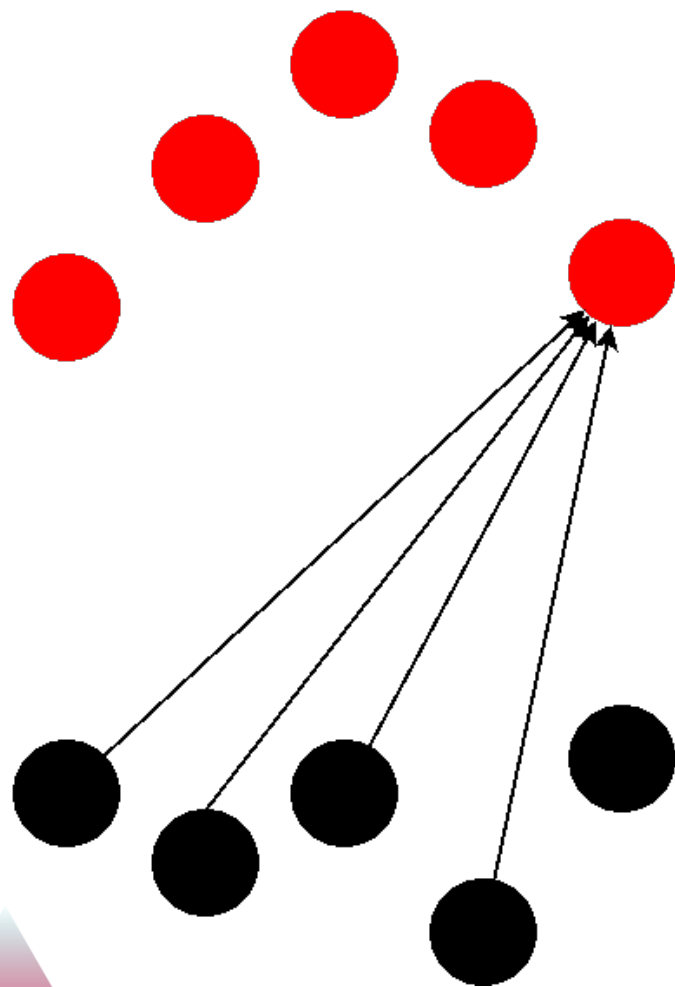
Can predict many
samples in parallel
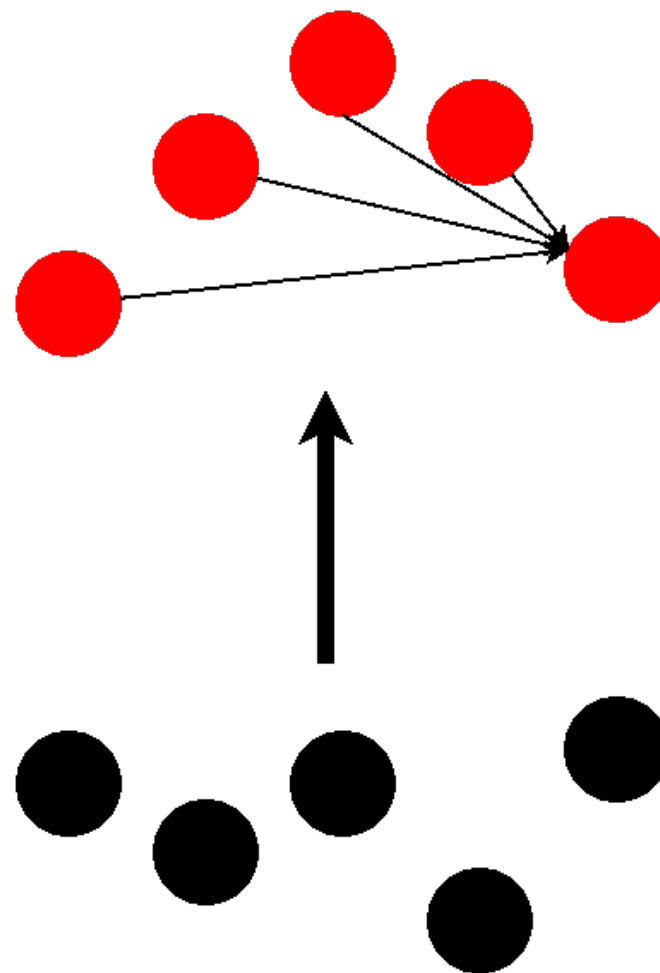
# Parallelized approach

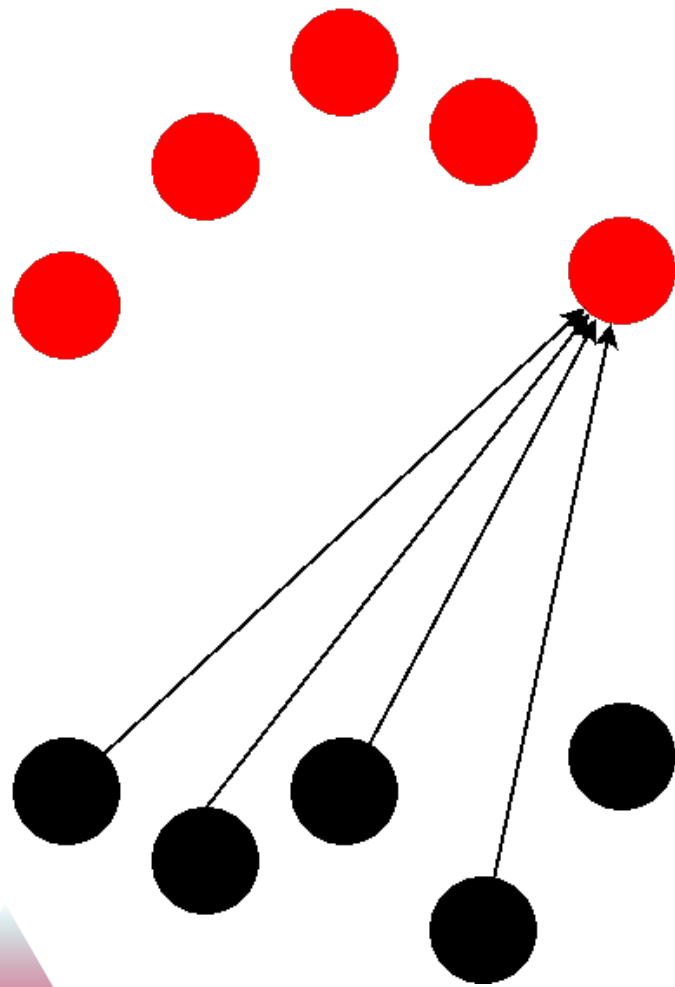Can predict many samples in parallel

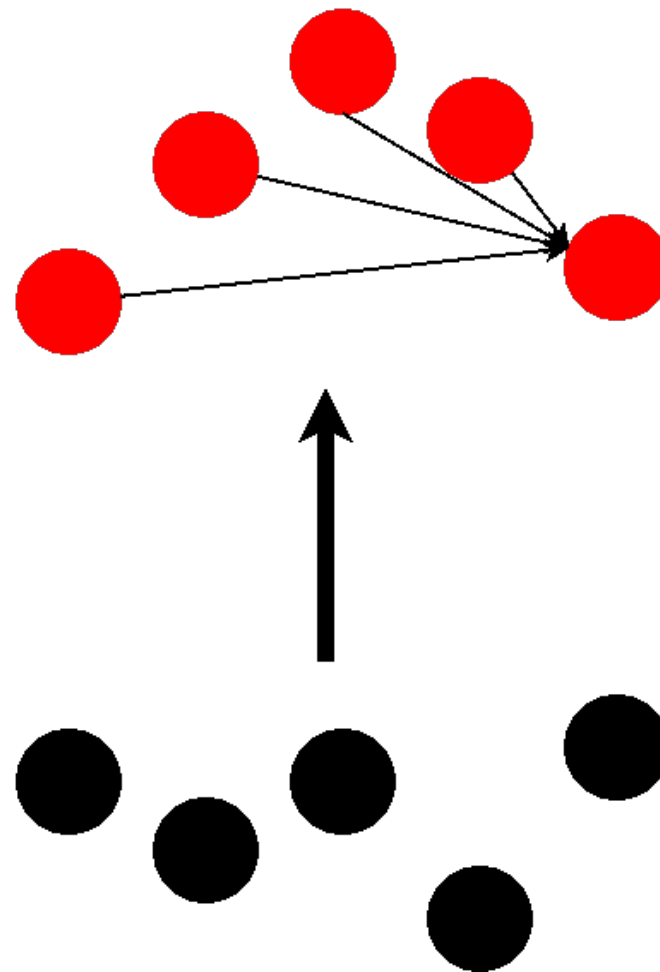Autoregressive model predicts likelihood of samples
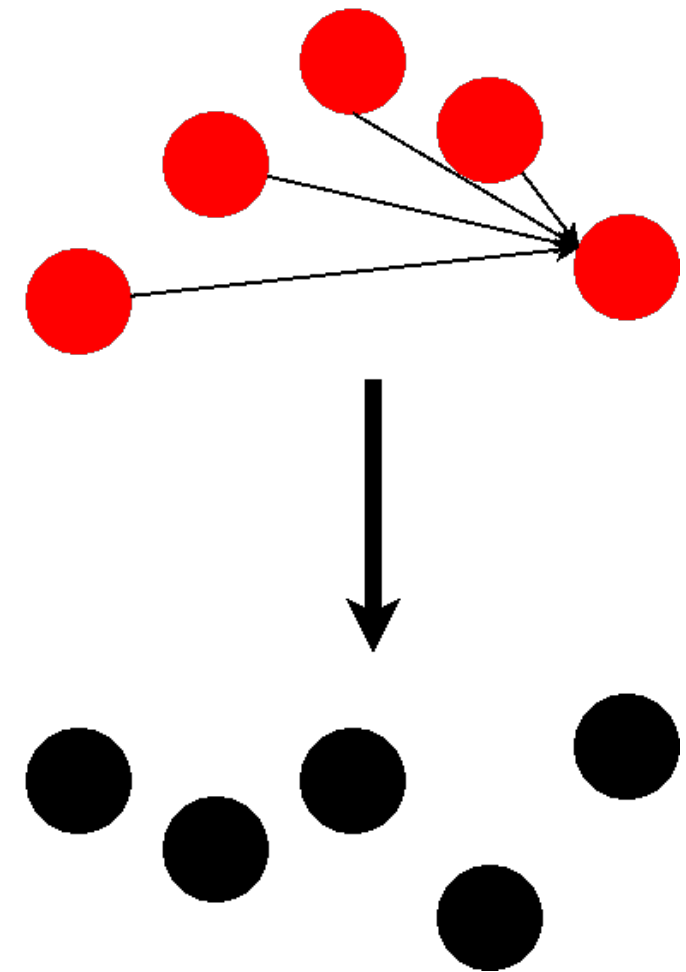


Reference: [b]

# Parallelized approach

Can predict many samples in parallel

Autoregressive model predicts likelihood of samples

Gradients maximizing this likelihood are passed back

# What do you do with this?  Voice Skins



Adversary

Generated Audio
targeting Sarena

Real Audio
from Sarena

# What do you do with this?  Voice Skins

Master Model

Speaker 1
Speaker 2
Speaker 3

Speaker 1
Speaker 2
Speaker 3

🔊 Input Speech

🔊 Output Speaker1

🔊 Output Speaker2

# Agenda

1. Overview of Neural Audio models
2. **Challenges for Development**
3. Challenges for Runtime
4. Epilogue

# Problem: Model training is time-consuming and costly



Training time of Speech Synthesis Models in Hours
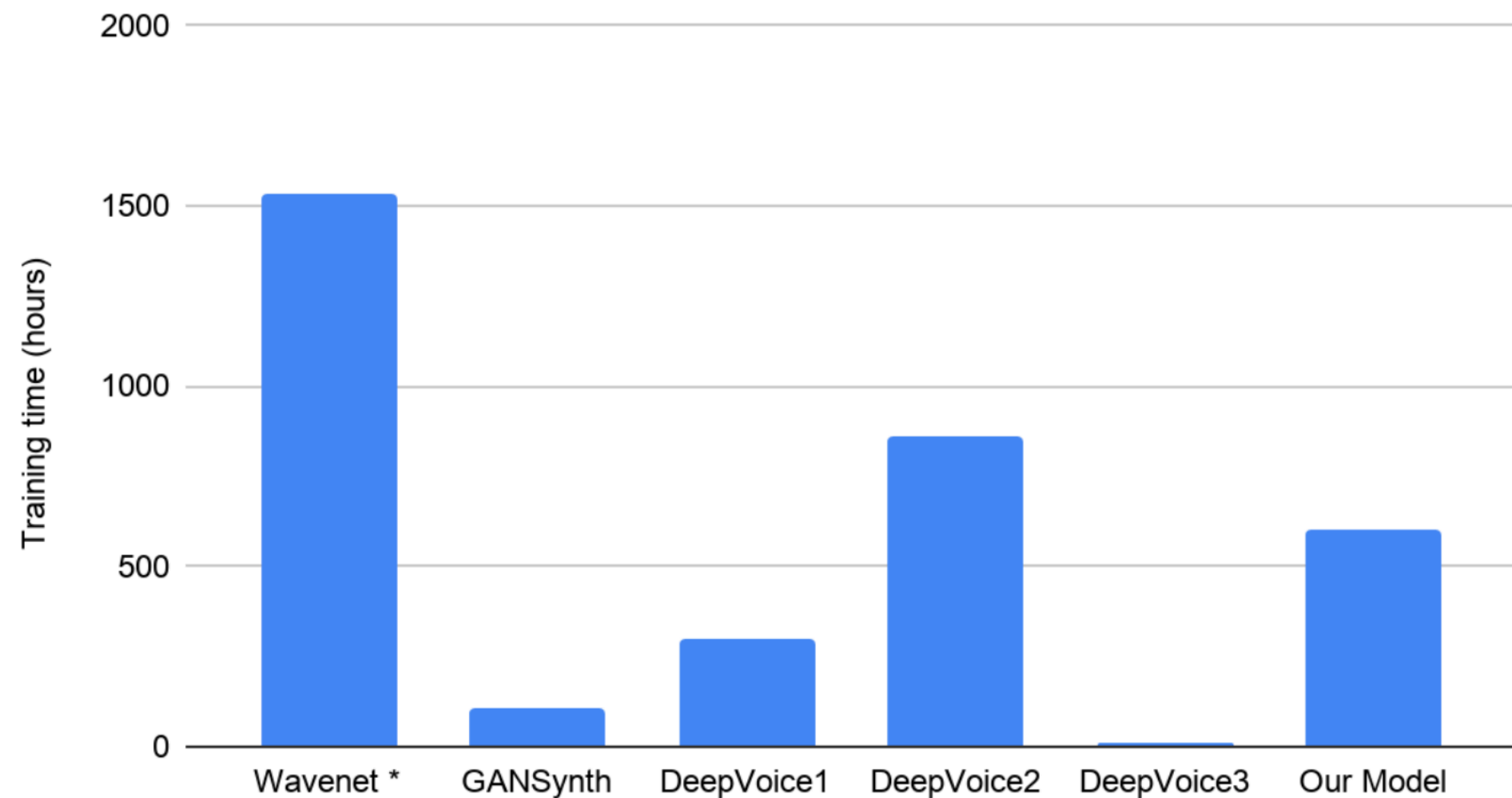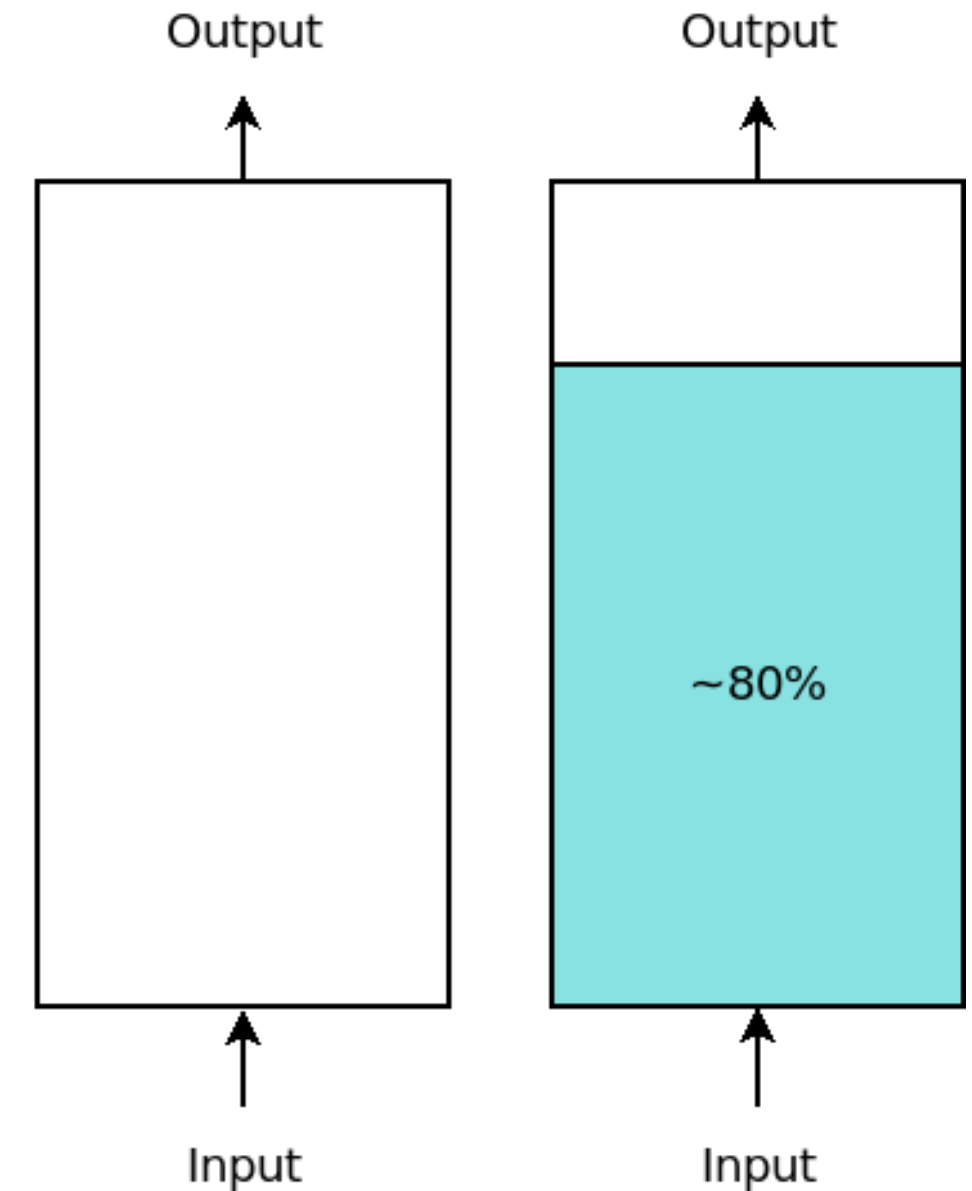
Reference: [c], [d], [e], [f], [g], [h]

# Solution: Train less and use sponsorship programs

## 1. Use pretrained models or stages of training!

- Multistage training on one v100 from 25 days to 5 days.
- Focuses experimentation on specific aspects of training

FOUNDER FRIENDLY LABS

### 2. Get Sponsorship!



Output          Output
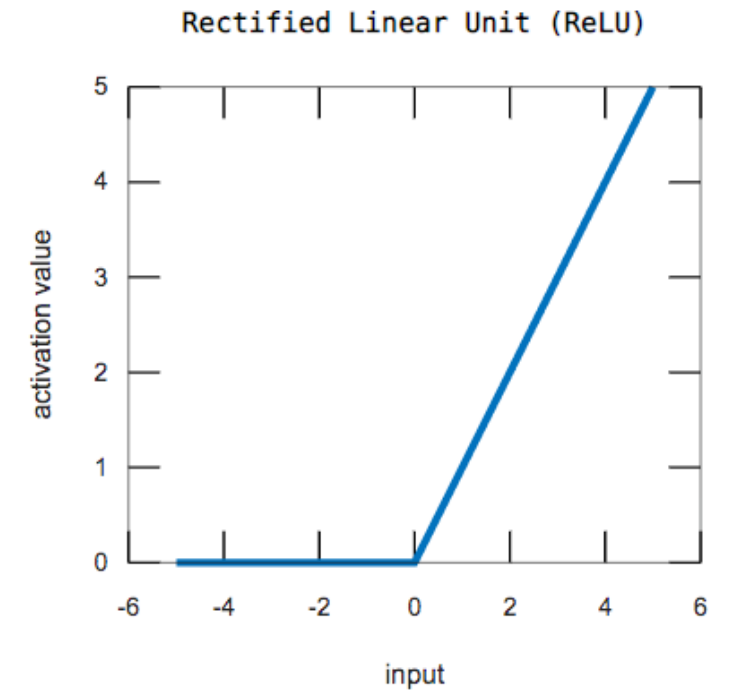
Frozen

Trained

~80%

Input          Input

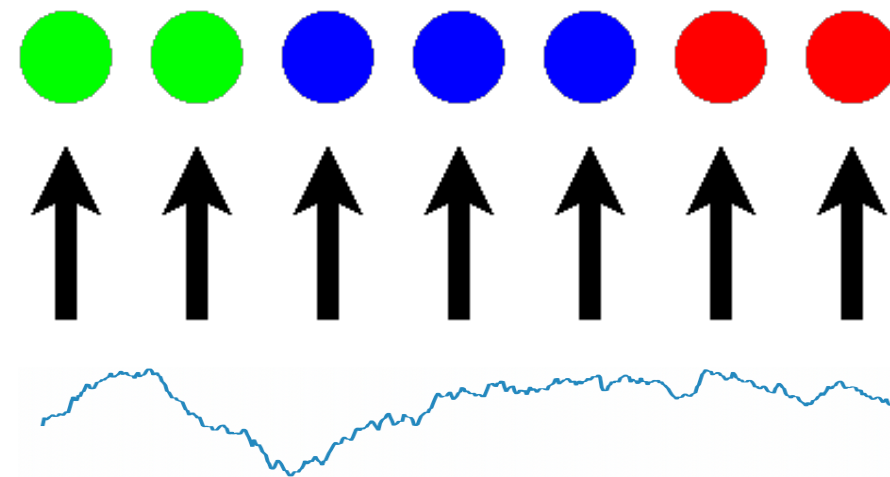# Solution: Kill experiments early based on heuristics

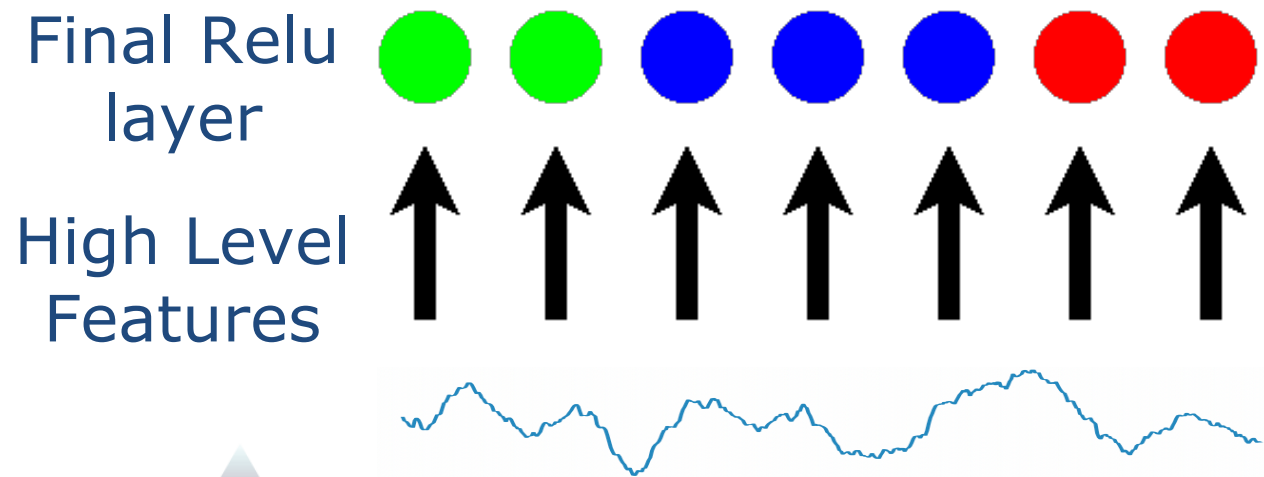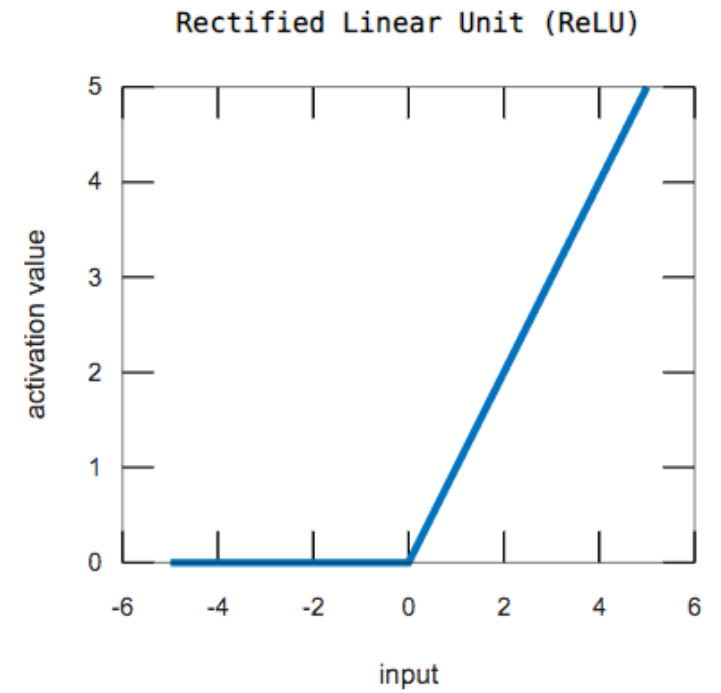Find correlations between negative results of experiments such as model collapse, poor generalization, poor performance

- Dead relu monitoring
- Other heuristics (power loss collapse, healthy average activation spread and more in [1])

# Example Solution: Dead Relus



Rectified Linear Unit (ReLU)

# Example Solution: Dead Relus

Rectified Linear Unit (ReLU)

Final Relu layer

High Level Features

# Example Solution: Dead Relus

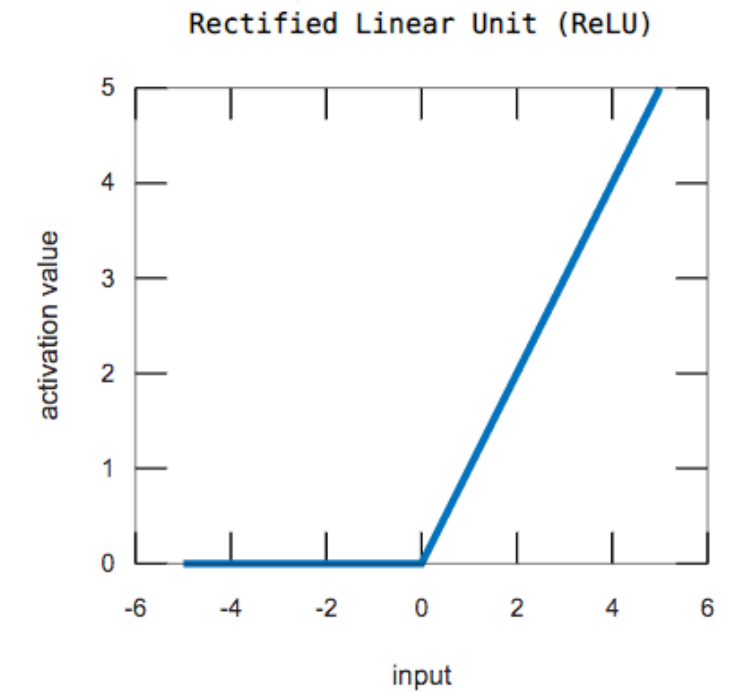🟢 Evergreen

🔵 Varies depending on input

🔴 Always 0 – Dead Relus

Rectified Linear Unit (ReLU)

Final Relu layer

High Level Features

# Example Solution: Dead Relus



Dead ReLus during Training

# Result of tracking Correlations

- Terminate failing experiments after hours or 1-2 days
- Iterate faster by spending less time heavily analyzing model performance

# Problem: How to Evaluate Models

# Problem: How to Evaluate Models



| Available Metrics | | |
|---|---|---|
| Training Losses | Other Neural Nets:<br>* Inception Score<br>* Frechet Distance | Human Score:<br>Mean Opinion Score (MOS) |
| - Difficult to weigh different losses, and doesn't always correlate with performance | - Highly dependent on the dataset and task for evaluation nets | - Slow/costly and difficult to target specific aspects of media |

# Solution: Design assessment tools for your needs

# Solution: Design assessment tools for your needs

Comparison of model behaviors along specific dimensions (e.g. input fidelity, intelligibility)

# Solution: Design assessment tools for your needs

Comparison of model behaviors along specific dimensions (e.g. input fidelity, intelligibility)

Simplifies collaboration

# Result: Sets of comparisons

- Takes ~15 minutes of time to compare an experiment with the baseline
- Streamlines team collaboration [2]

# Equipped with these tools, training models is possible!

1. Use pretrained models to shorten training time.
2. Use sponsorship if you're in a similar position to us.
3. Use heuristics to kill experiments early.
4. Design assessment tools for specific goals

# Agenda

1. Overview of Neural Audio models
2. Challenges for Development
3. **Challenges for Runtime**
4. Epilogue

# Why Do You Need to Run on Device

- **Low Latency Player Interaction**
  - **Audio feedback <= 30ms**
- Offline or Peer-to-Peer

# Why Do You Need to Run on Device

- Low Latency Player Interaction
  - Audio feedback <= 30ms
- **Offline or Peer-to-Peer**

# Constraints Running on Device

- **Reduce Resource Consumption**

- Different Device Capacities

- Optimizing for latency

- Audio Real-time friendly (~100Hz audio framerate)

Master Model Throughput

# Constraints Running on Device

- Reduce Resource Consumption
- **Different Device Capacities**
- Optimizing for latency
- Audio Real-time friendly (~100Hz audio framerate)

# Solution - Distillation

- Supervised training of smaller "distilled" model

# Solution - Distillation

- Reduce unneeded capacity

Distilled Model

Master Model

# Solution - Distillation

● Ship default & HQ versions of models

Throughput



10

8

6

4

2

0

| Master Model (1.6M params) | Distill HQ (430k params) | Distill Mid (320k params) | Distill Small (200k params) |

Realtime

Master Model

Distill Mid

(i7-6700K)

# Constraints Running on Device

- Reduce Resource Consumption
- Different Device Capacities
- **Optimizing for latency**
- Audio Real-time friendly (~100Hz audio framerate)

# Constraints Running on Device

- Reduce Resource Consumption
- Different Device Capacities
- Optimizing for latency
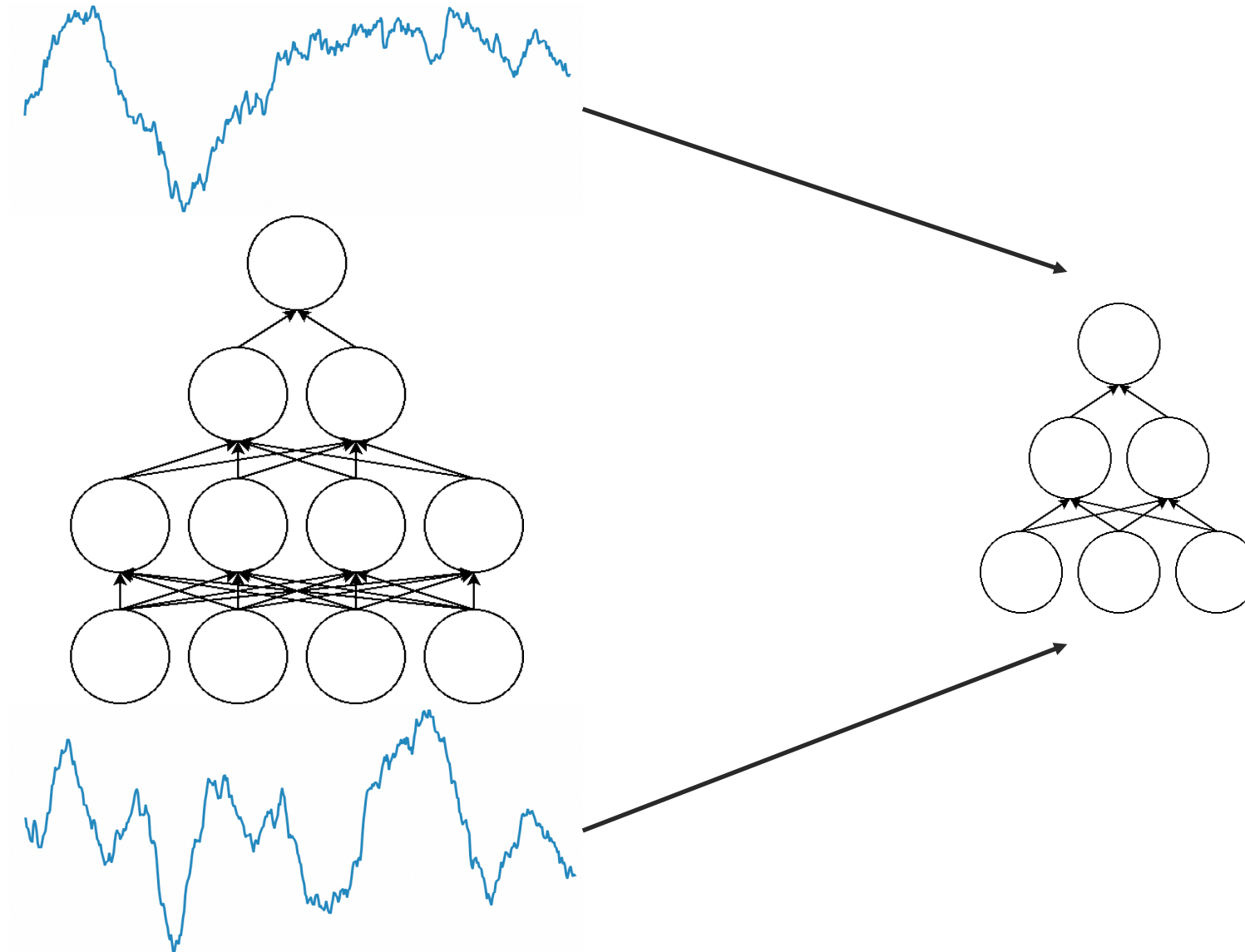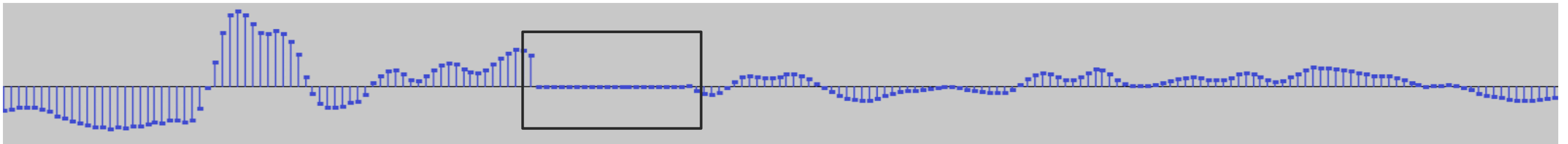- **Audio Real-time friendly (~100Hz audio framerate)**
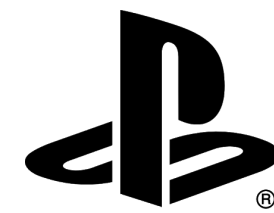
With Mallocs

# Solution - Custom Inference Code

- Custom C++ avoids memory allocations, blocking [3]
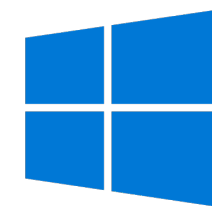- Use vectorization and fused-multiply-add

# Solution - Custom Inference Code

- Portable - cross-compiles to everything
- Easy integration, easy pipelines

iOS

macOS

PS4

Windows

REAPER
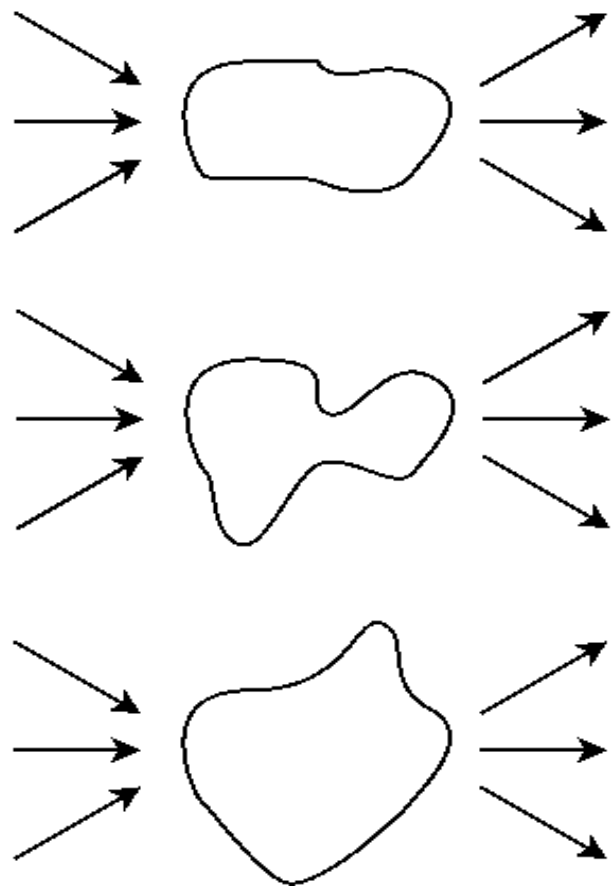Digital Audio Workstation

Wwise

NINTENDO SWITCH

XBOX

android

# Custom Code Time Investment

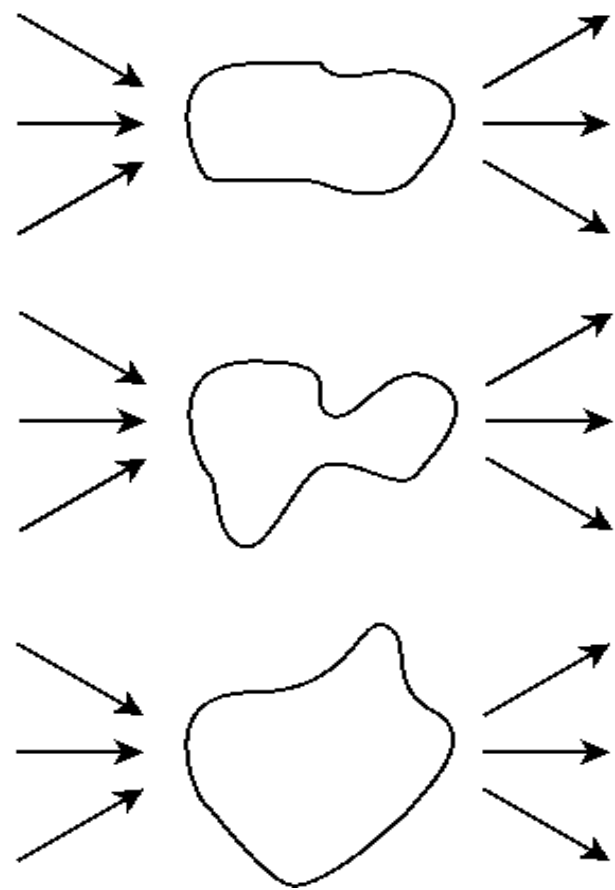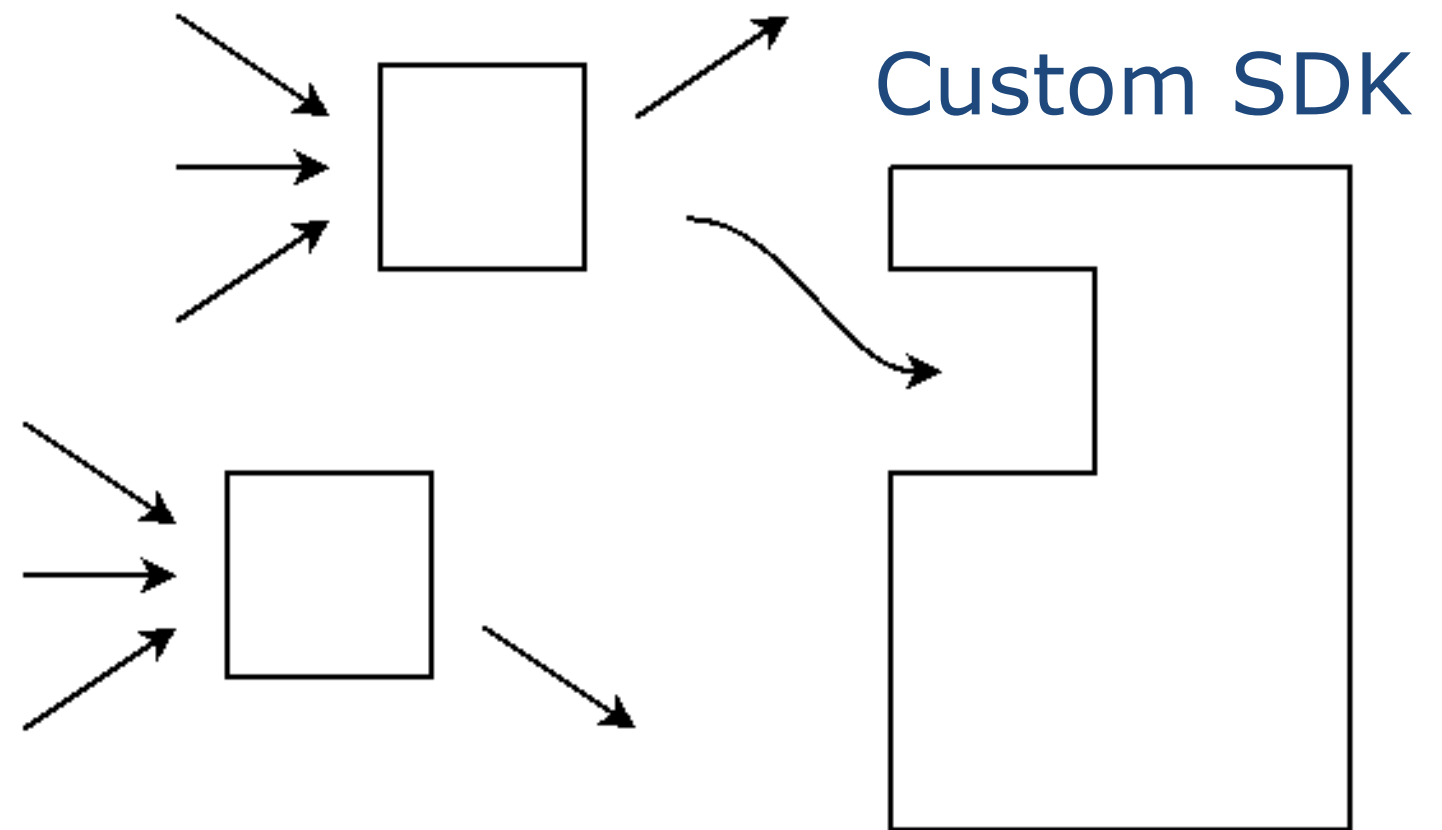- Distillation avoids rewrites [4]

Master Models

# Custom Code Time Investment

- Distillation avoids rewrites [4]
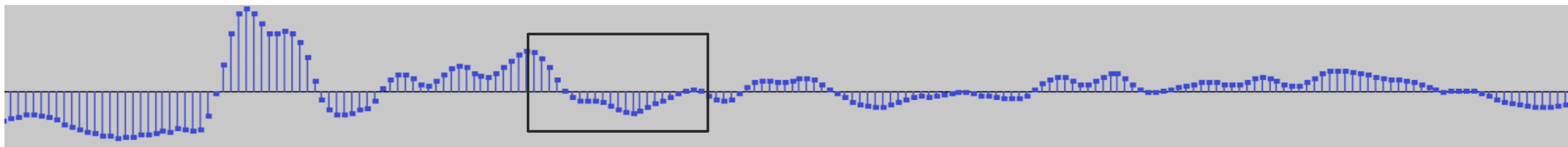
Master Models

Distilled Models

Custom SDK

# Custom Audio Code Tips

- Bundle in SRC
- Adapt to various frame sizes with flexible pre-allocation
- Include real-time audio logging [5]

# Custom Code Results

🔊 No Mallocs

# Agenda

1. Overview of Neural Audio models
2. Challenges for Development
3. Challenges for Runtime
4. **Epilogue**

# Summary

- Even a small team can train and deploy ML models
- Use Pretraining, Pruning, and QA tools
- Optimize with Distillation and Custom Inference Code

# Acknowledgements

- Also come talk to us for live demos!

# Referenced Posts

1. https://modulate.ai/blog-collection/2020/03/trainingheuristics

2. https://modulate.ai/blog-collection/2020/03/qatools

3. https://modulate.ai/blog-collection/2020/03/nonblockingconv

4. https://modulate.ai/blog-collection/2020/03/distillation

5. https://modulate.ai/blog-collection/2020/03/realtimeaudiologging

# Online References

a. https://www.deepmind.com/blog/article/wavenet-generative-model-raw-audio

b. https://arxiv.org/abs/1711.10433

c. https://cloud.google.com/compute/gpus-pricing

d. https://openreview.net/pdf?id=H1xQVn09FX

e. https://arxiv.org/pdf/1702.07825.pdf

f. https://aws.amazon.com/ec2/instance-types/p2/

g. https://arxiv.org/pdf/1705.08947.pdf

h. https://arxiv.org/abs/1710.07654