

Towards concurrent real-time audio-aware agents with deep reinforcement learning

Anonymous Full Paper
Submission 32

Abstract

Audio holds significant amount of information about our surroundings. It can be used to navigate, assess threats, communicate, as a source of curiosity, and to separate the sources of different sounds. Still, these rich properties of audio are not fully utilized by current video game agents.

We use spatial audio libraries in combination with deep reinforcement learning to allow agents to observe their surroundings and to navigate in their environment using audio cues. In general, game engines support rendering audio for one agent only. Using a hide-and-see scenario in our experimentation we show how support for multiple concurrent listeners can be used to parallelize the runtime operation and to enable using multiple agents. Further, we analyze the effects of audio environment complexity to demonstrate the scalability of our approach.

1 Introduction

Modern video games have rich and high-quality audio scenes. Such audio holds significant amount of information on the virtual environment. Deep learning is becoming increasingly popular as an approach to extract information from virtual environments, but is not fully utilized to its potential especially when audio is considered.

For video game non-player characters (NPCs), the utilization of audio could support intuitive, human-like behavior. For example, an audio-aware agent could have a fear of loud noises as a self preservation mechanic. Audio-aware agents could be curious about odd noises, gather information about the surrounding environment and become susceptible to audio-based distractions deliberately made by human players or even other NPCs.

While using deep learning for audio sensing in video games is not yet an extensively researched topic, there exists a wide range of applications that utilize it. For example, audio has been used to classify gunshot noises in video games [1], to enhance exploration during training [2] and as a generic addition to visual observations [3]. Games engines can integrate multi-disciplinary functionalities and have wide applicability beyond entertainment. As such, game engines can be seen as generic platforms [4] for research on deep learning methods.

However, game engines are often created with the assumption that audio is rendered for the human player only. As such, the environments can have only one audio listener at a time. This listener can either be assigned to the player or a single audio-aware agent. In order to support multiple concurrent audio-aware agents, advanced methods are needed.

In this paper, we describe our experimentation utilizing deep reinforcement learning (DRL). Our experimentation focuses on sound source localization and is based on a hide-and-see scenario, where agents must locate targets based on spatial audio cues. Our main contribution is to show that game engines can be enhanced by using multi-listener audio rendering to enable parallelization of the runtime operation and use multiple agents. We demonstrate how spatial audio libraries can be used in a popular game engine with machine learning support (the Unity engine) to construct a multi-listener based system that can be run in real-time on modern workstations and yield good performance even if the audio environments are complex. Our methodology is not limited to games. It can be used also in other domains that use similar virtual reality tools and have similar requirements. Our code is available in GitHub¹.

The structure of this paper is the following. We begin our presentation by reviewing related work in Section 2. In Section 3, we describe our methodology and our selection of tools. We continue by describing our experimental setup (Section 4) and the results from our experimentation (Section 5). We end our paper with a short discussion (Section 6) and our conclusions (Section 7).

2 Related work

Deep reinforcement learning (DRL) [5] is a branch of machine learning, where agents are trained in a trial-and-error manner. The training is done in an environment, where learning is based on the rewards the agents receive as feedback on their actions. As the agents are truly run during the training, making observations in the training environment and overall performance of the training steps are essential.

There exists a vast amount of machine learning literature on audio. Considering our application

¹<https://github.com/anonamee-333/anon-nldl-aaaa>

092 area, Grumiaux et al. [6] survey different deep learn- 145
 093 ing methods in general for sound source localization. 146
 094 The survey from Latif et al. [7] focuses only on the 147
 095 use of DRL but considers a wide range of audio- 148
 096 based applications. Beig et al. [8] review spatial 149
 097 sound rendering for virtual environments and games. 150

098 In applying DRL for game applications, timing 151
 099 behavior is central in addition to the related com- 152
 100 putational efficiency. Hedge et al. [3] study the 153
 101 performance of different model architectures in a 154
 102 series of tasks that require the agent to recognize 155
 103 sounds. They decouple audio rendering from dedi- 156
 104 cated sound hardware to enable faster-than-realtime 157
 105 parallel simulation in Vizdoom [9], resulting in faster 158
 106 training of audio-based agents. 159

107 Their work is an example of agents learning to play 160
 108 a game, instead of learning to act as NPCs. They 161
 109 reach a training throughput of 120 000 samples per 162
 110 second with audio and 150 000 samples per second 163
 111 without audio. 164

112 Cowan et al. [10] present a computationally effi- 165
 113 cient audio rendering system for game NPCs. They 166
 114 also recognize the problems caused by poor utiliza- 167
 115 tion of spatial audio by NPCs. Their main idea is 168
 116 that if the audio for the player is rendered using 169
 117 a graph-based sound propagation method, the gen- 170
 118 erated graph can immediately be reused for NPCs 171
 119 with low additional computational cost. 172

120 Environments including the interaction by the 173
 121 agents is typically modeled or simulated in a man- 174
 122 ner that captures the related physical realism. Chen 175
 123 et al. [11] present SoundSpaces 2.0, a simulator for 176
 124 training audio-based agents in environments mode- 177
 125 led after the real world. Gan et al. [12] show their 178
 126 results with discrete in-door navigation using audio- 179
 127 visual data with AI2-THOR platform. They create 180
 128 an offline dataset from AI2-THOR platform running 181
 129 on Unity game engine. 182

130 3 Audio agents in virtual envi- 183 131 ronments 184

132 In this section, we describe our methodology and the 185
 133 related selection of tools. We begin by discussing 186
 134 the problem space, and continue by describing our 187
 135 platform selection, the related observation tools, and 188
 136 the audio rendering libraries we have selected. After 189
 137 these, we describe our approach to multiple listeners, 190
 138 which are the key aspects of this research. 191

139 3.1 The problem space 192

140 While real-world embodied agents and video game 193
 141 agents can utilize audio using similar methods, they 194
 142 have a different set of requirements. 195

143 The audio-aware agents should be able to utilize 196
 144 audio in a useful way, but their output does not 197

145 need to be perfect. Unlike with real-world embod- 146
 147 ied agents, mistakes can even be beneficial for the 147
 148 gameplay experience. Additionally, the agents need 148
 149 to be able to work with soft real-time simulations. 149
 150 Missing a deadline is not critical, but can harm the 150
 151 experienced quality of service. The agents also do 151
 152 not need to make decisions on every simulation tick, 152
 153 but still often enough to be able to fluently navigate 153
 the simulated environment. 154

154 Considering the use of deep reinforcement learning 154
 155 for the audio-aware agents there are two challenges: 155
 156 1) The computational requirements of audio-aware 156
 157 agents and the related spatial audio rendering by 157
 158 the runtime system should be relatively low. Run- 158
 159 ning in real-time also during training avoids audio 159
 160 distortions. 2) The training of the agents should be 160
 161 efficient enough in order to make their use realistic 161
 162 in the context of video game development. Espe- 162
 163 cially the latter aspect is dependent on the operating 163
 164 environment of the agents as complex audio scenes 164
 165 make the training challenging. 165

166 3.2 Platform 166

167 We use Unity [4] as our platform. We combine Unity 167
 168 with its deep reinforcement learning framework ML- 168
 169 Agents [4]. ML-Agents includes a training pipeline 169
 170 for simple DRL architectures with popular methods, 170
 171 such as PPO [13]. It also exposes the Unity API for 171
 172 external tools, such as Ray RLib [14] enabling the 172
 173 use of complex DRL models. 173

174 3.3 Observations 174

175 Virtual agents observe their surroundings through 175
 176 abstractions called sensors. In the context on ML- 176
 177 Agents, the simplest sensors are fixed length vector 177
 178 sensors, that can be filled with arbitrary values. 178

179 To be able to use audio data as the input for 179
 180 our policy network, we need to have an audio sen- 180
 181 sor, but ML-Agents does not include any audio- 181
 182 related sensors. We opted to use a third-party sensor 182
 183 from GitHub user mbaske [15], which is available 183
 184 on GitHub with the MIT license. The sensor sup- 184
 185 ports observing stereo audio in spectral domain us- 185
 186 ing Short-Time Fourier Transform (STFT) [16] with 186
 187 different windowing functions [17], such as Hanning 187
 188 and Rectangular window. The windowing function 188
 189 affects the amplitude and frequency accuracy of the 189
 190 sensor. The sensor constructs an observation batch 190
 191 by sampling the audio into a buffer over several 191
 192 simulation ticks. 192

193 3.4 Audio rendering libraries 193

194 By default, popular game engines such as Unity 194
 195 and Unreal Engine support simplified spatialization 195
 196 without accounting for reflections or occlusions. To 196
 197 enable more complex audio observations, we use 197

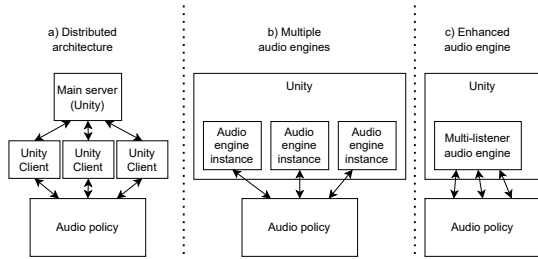


Figure 1. Approaches for multi-listener setup.

Table 1. System specifications

Label	CPU	GPU	RAM	OS
A	Intel 13900k	GTX 1080 ti	64 GB	Windows 11
B	Intel 13400	RTX 3080	32 GB	Ubuntu 22.04

198 Steam Audio [18] as a plugin for Unity. Steam Au-
 199 dio is one of the several libraries [8] that support
 200 spatial audio. It supports spatialization through
 201 head related transfer function (HRTF), which mod-
 202 els the interaural level differences (ILD) and the
 203 interaural time difference (ITD). It also supports
 204 occlusions and reflections either by real-time ray
 205 tracing for dynamic scenes or baked audio propaga-
 206 tion for static scenes.

207 3.5 Multi-listener support

208 Even with Steam Audio, Unity supports only one
 209 active audio listener at a time. This means that all
 210 audio will be rendered from the point of view of a
 211 single agent. From the training perspective, having
 212 only one audio listener limits parallelization of the
 213 training process as one instance will only support
 214 one agent. From the runtime perspective, the audio
 215 listener is often reserved for the player.

216 We propose three basic ways to enable multi-
 217 listener simulation. These are shown in Figure 1.

- 218 1. Each client has its own audio engine instance
 219 running and will therefore be able to feed mul-
 220 tiple agents with audio observations. The game
 221 logic runs on the main server, but the audio
 222 events are replicated on the remote clients.
- 223 2. Creating multiple instances of the audio engine
 224 within a single game instance, thus removing
 225 the overhead caused by the network.
- 226 3. An audio engine that natively supports multiple
 227 listeners, which eases internal optimizations.

228 Integrating the multi-listener support directly into
 229 the engines opens optimization opportunities and
 230 removes overhead but also increases significantly the
 231 implementation effort. We selected Option 1 for our
 232 experimentation.

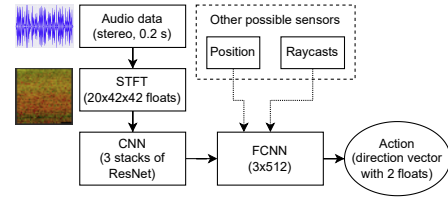


Figure 2. Sensor data processing in our agent design.

4 Experiment setup

233

234 In this section, we describe the setup for our ex-
 235 periment. Our goal is to create and evaluate the
 236 effect of using multiple listeners as proposed in Fig-
 237 ure 1 (leftmost), while also training and evaluating
 238 a DRL-based audio-aware agent using Unity and
 239 Steam Audio. The multi-listener measurements are
 240 done on hardware configuration A as seen in Table
 241 1. The agents are trained using both configurations
 242 A and B.

243 We consider a hide-and-seek scenario, where an
 244 audio-aware agent attempts to find the target (e.g.,
 245 the player) in an indoors area. The target is con-
 246 stantly making noise by randomly playing back one
 247 of eight footstep audio clips. In order to make sim-
 248 plify the evaluation of the agents, the target is rep-
 249 resented as a static object that does not move around
 250 in the environment. The motivation for this setup
 251 comes from a typical scenario, where the player at-
 252 tempts to evade detection.

253 The agent can observe its surroundings with sen-
 254 sors as described in Figure 2. The audio data is fed
 255 through STFT to a convolutional neural network
 256 (CNN) and the integration of other sensor data is
 257 done with a fully connected neural network (FCNN).
 258 In our experiment, the other sensors are excluded
 259 and the audio sensor is the only source of data.

4.1 Multi-listener setup

260

261 We measure the impact of running multiple simu-
 262 lation instances in parallel as proposed in Figure
 263 1 (leftmost). Each Unity instance is running the
 264 exact same scene with Steam Audio and one audio
 265 agent with a CNN-based audio model. The number
 266 of audio sources is varied between 1, 10 and 30 au-
 267 dio sources to evaluate the impact of having several
 268 audio sources.

4.2 Environments

269

270 The agent performance is compared in multiple en-
 271 vironments with increasing complexity. The easiest
 272 environment ("Simple") is an empty hall without any
 273 obstructions. The second environment ("Medium")
 274 adds some obstructions in the form of rooms and
 275 walls, making it necessary to utilize indirect audio

276 reflections. The third environment ("Complex") in-
277 creases the difficulty by including more walls. The
278 environments are described in more detail in Ap-
279 pendix B.

280 4.3 Agent design

281 The agents are trained using PPO [13] over 10 mil-
282 lion samples in the second ("Medium") environment.
283 The training hyperparameters are listed in Appendix
284 A.1.

285 The policy design is shown in Figure 2. In order to
286 isolate the usefulness of the audio observations, the
287 audio sensor is the only sensor included in the experi-
288 ments. The agent outputs a two values, representing
289 the x and y values of a 2-dimensional vector.

290 The reward signal, denoted as R , is calculated by
291 taking a dot product between the direction of the
292 shortest path to the target and the action vector. To
293 encourage finishing episodes as quickly as possible,
294 we convert the result into a penalty by shifting it
295 from range $[-1, 1]$ to range $[-2, 0]$. Finally, it is scaled
296 by the fixed delta time to make it invariant to the
297 simulation physics update frequency. The agent is
298 given an additional reward of 2 upon finishing an
299 episode.

300 The reward R at physics update n can be ex-
301 pressed as:

$$302 R_n = ((\mathbf{a} \cdot \hat{\mathbf{s}}) - 1)\Delta t + \begin{cases} 2, & \text{if target reached} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

303 where \mathbf{a} is the action vector, ($\hat{\mathbf{s}}$) is a unit vector
304 pointing to the shortest path to target and Δt is
305 the time between two physics updates (called fixed
306 delta time in Unity). Maximum achievable reward
307 during an episode is 2.

308 To prevent agents from getting stuck between two
309 locations, they need information on locations already
310 visited. Instead of using more complex neural net-
311 works like RNNs [19] or LSTMs, we used a simple
312 weighted navigation grid (Appendix A.3). The grid
313 node weights indicate the likelihood of audio coming
314 from that direction. We adjust these weights based
315 on neural network outputs and navigate towards the
316 highest-weighted neighboring node. If there are no
317 occlusions between the target and the agent, the
318 navigation grid is ignored and the agent will travel
319 directly towards the output direction.

320 4.4 Evaluation metrics

321 We evaluate the viability of our multi-listener ap-
322 proach by measuring the CPU and memory (MEM)
323 utilization with different configurations.

324 We measure the agent performance using sample
325 throughput, training time, and cumulative training
326 reward during training. After training, we evaluate

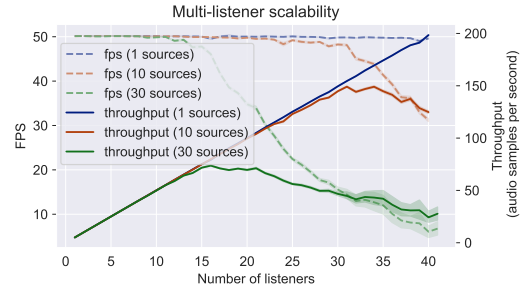


Figure 3. Main server FPS and total audio sample throughput. Each audio sample contains a 0.2 second audio buffer collected over 10 simulation ticks.

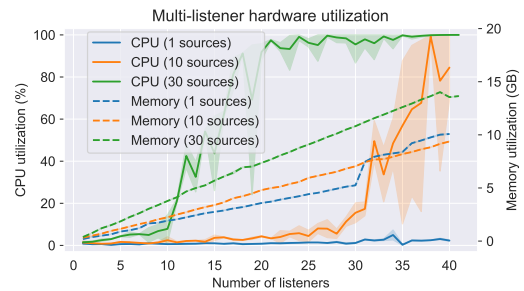


Figure 4. Average hardware utilization. CPU utilization is an average over all 24 cores. Both metrics are also affected by background processes.

327 the agent with a fixed set of 100 randomized episodes
328 in each environment using the SPL metric [20].

329 SPL is defined as

$$330 \text{SPL} = \frac{1}{N} \sum_{i=1}^N S_i \frac{l_i}{\max(p_i, l_i)} \quad (2)$$

331 where N is the number of episodes, S_i is 1 if the
332 episode i succeeded and 0 otherwise. The l_i is the
333 shortest path to target and p_i is the length of the
334 path taken by the agent.

335 As such, SPL accounts for both failure to reach
336 the goal and the optimality of the path compared to
337 the shortest path. It is better than simply measuring
338 number of steps or time to reach the target, since
339 the distance to the target varies by episode. However,
340 as SPL relies on knowing the optimal path length,
341 it is not suitable for moving targets as such.

342 5 Results

343 In this section, we describe our measurement results
344 on multi-listener performance, training performance,
345 and agent performance with some analysis on the
346 probable causes of the observed behavior.



Figure 5. Training reward over time.

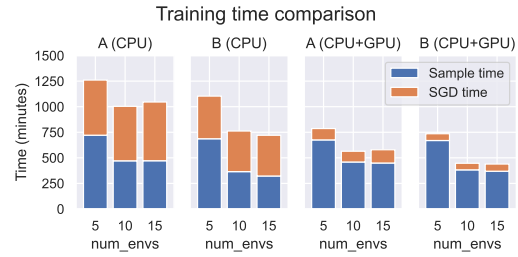


Figure 6. Training time comparison with different computation configurations.

347 5.1 Multi-listener performance

348 Figure 3 shows how the performance scales with
349 parallel audio-aware agents in three different config-
350 urations. The figure shows the simulation frames per
351 second (FPS) processed by the main server, which
352 is set at a fixed rate of 50 FPS. Values lower than
353 50 means that the simulator can no longer keep up
354 with real-time. With only one audio source, the FPS
355 stays nearly stable at up to 40 listeners. With 30
356 sources, the simulation supports roughly 12 concur-
357 rent audio-aware agents.

358 Figure 3 also shows the total sample throughput
359 for all audio-aware agents. Each audio sample is
360 computed from a 0.2 second long audio buffer, which
361 requires 10 simulation frames at 50 FPS. Sample
362 throughput often directly correlates with the time it
363 takes to train a DRL-based agent. From the figure,
364 we can see that the sample throughput increases
365 linearly to 200 audio samples per second with one
366 audio source. With 30 audio sources, the throughput
367 caps at 75 at around 15 environments.

368 Figure 4 shows how the utilization of the CPU
369 scales with parallel Unity instances. With one audio
370 source, the CPU utilization stays below 5 % even
371 at 40 parallel instances. With 10 audio sources,
372 the CPU utilization stays reasonably low until 25
373 environments before increasing exponentially. With
374 30 audio sources, the CPU usage increases already
375 after 10 instances. At 41 parallel instances, the
376 whole system becomes unresponsive.

377 Figure 4 also shows that each instance has a near
378 linear effect on memory usage, which is on par with
379 expectations. However, the number of audio sources
380 has a clear effect on the coefficient of this linear
381 behavior. With 30 audio sources the setup requires
382 nearly double the memory compared to having one
383 source. With 30 audio sources, one instances con-
384 sumes roughly 300 MB of memory. The sudden
385 increase of memory usage with one audio source
386 after 30 listeners is likely caused by background load
387 related to the Windows environment.

388 5.2 Training performance

389 We trained two DRL-based agents using two different
390 STFT window functions (STFT-R for rectangular

and STFT-H for Hanning window). The training 391
was repeated three times. During training, we collect 392
a sample on every simulation step. This increases 393
the sampling throughput, but the samples will have 394
overlap with each other as the audio buffer holds 395
data over 10 steps. 396

The Figure 5 shows the training reward curves. 397
The curves are averaged over all three runs and then 398
smoothed. The reward curves converge fast and 399
there is little improvement after 5 million steps. 400

Figure 6 shows the total training time for 10 mil- 401
lion steps with four different configurations. In ad- 402
dition to the total time, it shows how much time is 403
spent in collecting samples from the agents and how 404
long is spent updating the model with the Stochastic 405
Gradient Descent (SGD) [13]. With these config- 406
urations, increasing the amount of parallel Unity 407
instances up to 10 significantly increases the sample 408
throughput. Additionally, using GPU can nearly 409
halve the total training time by drastically reducing 410
time spent in computing SGDs. 411

Despite having a higher-end CPU, the configura- 412
tion A (CPU) is slower than B (CPU). This could 413
be due to the different operating system or some 414
other differences in hardware details, such as the 415
memory speed or motherboard configuration. 416

In conclusion, the sample throughput is a clear 417
bottleneck in a GPU-accelerated system. It might 418
be possible to further increase the sample through- 419
put by having multiple agents in one Unity instance 420
and by having faster-than-realtime audio render- 421
ing. Both of these could reduce the overhead from 422
running multiple Unity instances in parallel. 423

5.3 Audio-aware agent performance 424

425 After training the STFT-R and STFT-H agents, we
426 compared them to a random agent that traverses
427 randomly across the environment. As we repeated
428 the training three times, the results are also averaged
429 over three evaluations.

430 The Figure 7 shows a comparison of SPL values
431 achieved by different agents. The best possible SPL
432 value is 1. As expected, the audio-aware agents per-
433 form best in the simple scene with no obstacles with
434 an approximate SPL value of 0.85. The STFT-H

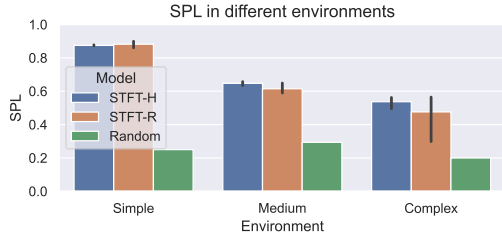


Figure 7. Measured SPL in different testing environments. Higher values are better.

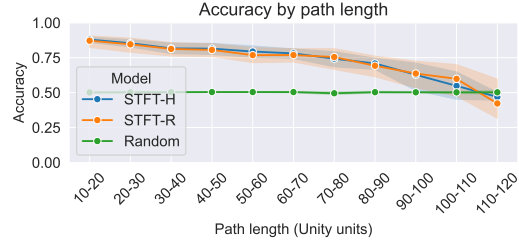


Figure 9. Accuracy by remaining shortest path length.

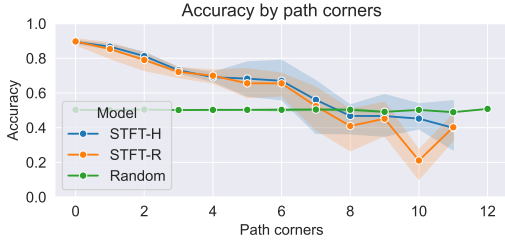


Figure 8. Accuracy by remaining shortest path corners.

seems to perform slightly better than SRFT-R overall. When obstacles are added, the agents perform worse, but still significantly better than the random agent. As reference, Anderson et al. [20] mention that an SPL of 0.5 shows a good level of navigation performance in reasonably complex previously unseen environments.

Figures 8 and 9 show how the remaining distance from the agent to the target affects agent prediction accuracy. If we denote the accuracy as $A \in [0, 1]$, then $A = 1$ means that the action-vector is aligned with the shortest path to target. Likewise, $A = 0$ means that the action-vector deviates from the direction of the shortest path by 180 degrees.

Figure 8 shows the accuracy based on how many corners the shortest path to target has. The number of corners depends on how many obstacles are between the agent and the target, since the agent must turn at least twice to get around an obstacle. Figure 9 shows the length of the shortest path to target in Unity units. From these figures, we can see that the remaining distance has a more linear effect on accuracy, while the amount of remaining obstacles has a more drastic effect.

6 Discussion

The lack of support for multiple audio listeners is a significant drawback in most current game engines as it limits the options of using audio as an observation for deep learning agents.

However, it seems that our approach of using multiple clients can ease these limitations. Our experimentation demonstrates that with a rather simple

workaround in Unity nearly a dozen of audio-aware agents can be supported. However, resource intensive games likely need more game-specific optimization than what we have done in our experimentation.

Our results also highlight that minor details, such as the STFT windowing function, can noticeably affect agent behavior. Therefore, the audio properties of the environment and the audio sensor must be chosen carefully.

As future work, we suggest looking at the other proposed approaches for creating multi-listener support to reduce overhead. Having only one game instance would remove the need for synchronization between the clients and reduce memory usage. Similar to Cowan et al. [10] it is also likely that some of the audio propagation computation could be reused.

Additionally, modifying the audio engine to support faster-than-realtime playback would be helpful. For example, in our experiment, the sample throughput was around 250 samples per second. This is rather slow compared to the 120 000 samples per second reached with the faster-than-realtime solution by Hegde et al. [3] for VizDoom. Similar to the solution by Hedge et al., the audio backend would likely need to be modified to enable full software rendering of the audio at any sampling rate

7 Conclusion

In this paper, we presented our experimentation on using multi-listener rendering in game engines for enabling concurrent audio-aware agents with deep reinforcement learning. Our work demonstrates how current audio libraries can be utilized to create audio-aware agents with reasonable training results.

While our results show the viability of using deep learning based audio-aware agents in video games, our research also highlights the need for further research on the methodology and development of applications and tools. The ability of faster-than-realtime audio operations could significantly speed up training of audio-aware agents. Likewise, native support for multiple audio listeners would make training and runtime more efficient while enabling the use of multiple agents.

References

- 510
- 511 [1] J. Park, Y. Cho, G. Sim, H. Lee, and J. Choo.
512 “Enemy spotted: In-game gun sound dataset
513 for gunshot classification and localization”. In:
514 *2022 IEEE Conference on Games (CoG)*. Bei-
515 jing, China: IEEE, Aug. 2022. DOI: [10.1109/
516 cog51982.2022.9893670](https://doi.org/10.1109/cog51982.2022.9893670).
- 517 [2] C. Gan, X. Chen, P. Isola, A. Torralba,
518 and J. B. Tenenbaum. “Noisy Agents: Self-
519 supervised Exploration by Predicting Audi-
520 tory Events”. In: *2022 IEEE/RSJ Interna-
521 tional Conference on Intelligent Robots and
522 Systems (IROS)*. Oct. 2022, pp. 9259–9265.
523 DOI: [10.1109/IROS47612.2022.9981614](https://doi.org/10.1109/IROS47612.2022.9981614).
- 524 [3] S. Hegde, A. Kanervisto, and A. Petrenko.
525 “Agents that listen: High-throughput reinforc-
526 ement learning with multiple sensory systems”.
527 In: *2021 IEEE Conference on Games (CoG)*.
528 IEEE, Aug. 2021. DOI: [10.1109/cog52621.
529 2021.9619096](https://doi.org/10.1109/cog52621.2021.9619096).
- 530 [4] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J.
531 Harper, C. Elion, C. Goy, Y. Gao, H. Henry, M.
532 Mattar, and D. Lange. “Unity: A general plat-
533 form for intelligent agents”. In: *arXiv* (2018).
534 arXiv: [1809.02627](https://arxiv.org/abs/1809.02627) [cs.LG].
- 535 [5] R. S. Sutton and A. G. Barto. *Reinforcement*
536 *Learning: An Introduction*. en. MIT Press, Nov.
537 2018. ISBN: 9780262039246.
- 538 [6] P.-A. Grumiaux, S. Kitić, L. Girin, and A.
539 Guérin. “A survey of sound source localization
540 with deep learning methods”. en. In: *J. Acoust.*
541 *Soc. Am.* 152.1 (July 2022), p. 107. ISSN: 0001-
542 4966, 1520-8524. DOI: [10.1121/10.0011809](https://doi.org/10.1121/10.0011809).
- 543 [7] S. Latif, H. Cuayáhuítl, F. Pervez, F.
544 Shamshad, H. S. Ali, and E. Cambria. “A
545 survey on deep reinforcement learning for
546 audio-based applications”. In: *Artificial In-
547 telligence Review* 56.3 (Mar. 2023), pp. 2193–
548 2240. ISSN: 1573-7462. DOI: [10.1007/s10462-
549 022-10224-2](https://doi.org/10.1007/s10462-022-10224-2).
- 550 [8] M. Beig, B. Kapralos, K. Collins, and P.
551 Mirza-Babaei. “An Introduction to Spatial
552 Sound Rendering in Virtual Environments and
553 Games”. In: *The Computer Games Journal*
554 8.3 (Dec. 2019), pp. 199–214. ISSN: 2052-773X.
555 DOI: [10.1007/s40869-019-00086-0](https://doi.org/10.1007/s40869-019-00086-0).
- 556 [9] M. Kempka, M. Wydmuch, G. Runc, J.
557 Toczek, and W. Jaśkowski. “ViZDoom: A
558 Doom-based AI Research Platform for Vi-
559 sual Reinforcement Learning”. In: *IEEE Con-
560 ference on Computational Intelligence and
561 Games*. Santorini, Greece: IEEE, Sept. 2016,
562 pp. 341–348. DOI: [10.1109/CIG.2016.
563 7860433](https://doi.org/10.1109/CIG.2016.7860433).
- [10] B. Cowan, B. Kapralos, and K. Collins. “Re-
564 alistic Audio AI: Spatial Sound Modelling to
565 provide NPCs with Sound Perception”. In:
566 *Audio Engineering Society Audio for AR/VR*.
567 Aug. 2020. URL: [https://aes2.org/
568 publications/elibrary-page/?id=20876](https://aes2.org/publications/elibrary-page/?id=20876).
569
- [11] C. Chen, C. Schissler, S. Garg, P. Kobernik, A.
570 Clegg, P. Calamia, D. Batra, P. W. Robinson,
571 and K. Grauman. “SoundSpaces 2.0: A Simu-
572 lation Platform for Visual-Acoustic Learning”.
573 In: *arXiv* (2022). arXiv: [2206.08312](https://arxiv.org/abs/2206.08312) [cs.SD].
574
- [12] C. Gan, Y. Zhang, J. Wu, B. Gong, and
575 J. B. Tenenbaum. “Look, Listen, and Act:
576 Towards Audio-Visual Embodied Navigation”.
577 In: *arXiv* (2019). arXiv: [1912.11684](https://arxiv.org/abs/1912.11684) [cs.CV].
578
- [13] J. Schulman, F. Wolski, P. Dhariwal, A. Rad-
579 ford, and O. Klimov. “Proximal Policy Opti-
580 mization Algorithms”. In: *arXiv* (2017). arXiv:
581 [1707.06347](https://arxiv.org/abs/1707.06347) [cs.LG].
582
- [14] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R.
583 Fox, K. Goldberg, J. Gonzalez, M. Jordan, and
584 I. Stoica. “RLlib: Abstractions for distributed
585 reinforcement learning”. In: *International con-
586 ference on machine learning*. PMLR. 2018,
587 pp. 3053–3062.
588
- [15] mbaske. *Audio Sensor Component for Unity*
589 *ML-Agents*. en. URL: [https://github.com/
590 mbaske/ml-audio-sensor](https://github.com/mbaske/ml-audio-sensor).
591
- [16] B. Boashash. *Time-frequency signal analysis*
592 *and processing: A comprehensive reference*. en.
593 2nd ed. San Diego, CA: Academic Press, Dec.
594 2015. ISBN: 9780123984999. DOI: [10.1016/
595 c2012-0-00024-5](https://doi.org/10.1016/c2012-0-00024-5).
596
- [17] O. Özhan. “Short-Time-Fourier Transform”.
597 en. In: *Basic Transforms for Electrical*
598 *Engineering*. Cham: Springer Interna-
599 tional Publishing, 2022, pp. 441–464.
600 ISBN: 9783030988456,9783030988463. DOI:
601 [10.1007/978-3-030-98846-3_7](https://doi.org/10.1007/978-3-030-98846-3_7).
602
- [18] Valve. *Steam Audio*. [https://
603 valvesoftware.github.io/steam-audio/](https://valvesoftware.github.io/steam-audio/).
604 Accessed: 2024-9-5.
605
- [19] I. Goodfellow, Y. Bengio, and A. Courville.
606 *Deep Learning*. en. MIT Press, Nov. 2016. ISBN:
607 9780262035613.
608
- [20] P. Anderson, A. Chang, D. S. Chaplot, A.
609 Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka,
610 J. Malik, R. Mottaghi, M. Savva, and A. R. Zamir.
611 “On Evaluation of Embodied Navigation
612 Agents”. In: *arXiv* (2018). arXiv: [1807.06757](https://arxiv.org/abs/1807.06757)
613 [cs.AI].
614
- [21] Valve. *Steam Audio Material documentation*.
615 [https://valvesoftware.github.io/
616 steam-audio/doc/unity/material.html](https://valvesoftware.github.io/steam-audio/doc/unity/material.html).
617 Accessed: 2024-9-6.
618

A Agent details

A.1 Training hyperparameters

The agents are trained using Unity ML-Agents library with PPO. The structure of the neural network (NN) follows the default implementation included in ML-Agents (version 0.30.0) with hyperparameter setting given in A.1. The audio data is included as a set of images, which ML-Agents feeds to a Resnet-style CNN. The output of the CNN is concatenated with the vector inputs from other sensors and fed to a linear network. In the experiments presented in this paper, the other sensors were excluded and the agents were using the audio sensor only.

Parameter	Value
maximum steps	10 000 000
batch size	1024
buffer size	10240
learning rate	0.0003
beta	0.005
epsilon	0.2
lambda	0.90
num epoch	3
hidden units	512
num layers	3

Table A.1. ML-Agents hyperparameters

A.2 Random agent

The random agent does not have any NN-based model. Instead, it chooses random coordinates in the environment and navigates towards them using the Unity NavMesh. Upon reaching these coordinates, it selects new coordinates. This process is repeated until the agent gets within 15 units of the target.

When the agent arrives within 15 units of the target, we assume that it can see the target and the agent navigates directly towards the target. Without having this kind of an detection range, it would be unlikely for the agent to ever encounter the target. The environments are 100x100 units large.

A.3 Navigation grid

Without any assistance, the DRL-based agents can easily get stuck between two points. Typically, this happens near walls: an agent infers that the audio is coming from the direction of the wall. The agent gets stuck in a loop of retreating and approaching. The behavior can be corrected by keeping a record of historical actions and observations, which we implement with a navigation grid (**Navgrid**).

Navgrid is a grid of evenly spaced navigation nodes placed across the environment. Adjacent nodes are neighbours to each other, if there is a clear unobstructed in line of sight between the nodes. Instead of directly using the NN output for navigation, we use it to update the **Navgrid** weights and navigate towards the highest-weight neighboring node. When a node is reached, the agent starts navigating towards the neighboring node that leads to the path of highest cumulative discounted weights.

Each node has a weight $W \in [0, 100]$. To represent the uncertainty of old weights, all node weights decay towards a neutral value $W_{neutral} = 20$ at the rate of $\gamma = 0.1$ units per physics update.

The pseudocode for the weight update algorithm is visible in Listing 1. The algorithm updates the weights of nearby nodes at the maximum depth of $d_{max} = 2$ nodes. Nodes in the general direction

Listing 1. NavGrid update pseudocode

```
def UpdateWeights(agent, action):
    c = agent.getClosestNode()
    c.setWeight(0)
    N = c.neighbors
    for n in N:
        dir = agent.getUnitVectorToNode(n)
        w = dot(action, dir) # [-1, 1]
        n.addWeight(w)
    U = n.neighbors
    for u in U:
        if u in N:
            continue
        u.addWeight(w * discount)
```

of the action will get an increase in weight, while nodes in the opposite direction will get a reduction in weight. The change in weight is propagated to further neighbors with a discount of λ^{d-1} (where d is the depth and $\lambda = 0.5$) up to the maximum depth d_{max} , as the uncertainty of the audio direction increases the deeper we progress into the graph.

The weight of the closest node to the agent is set to 0, as we can be certain that the target is not nearby. This is because the agent does not use the **Navgrid**, if it has line-of-sight to the target. In our experiment, the agent is given the line-of-sight information, but it could also be predicted by the agent from sensor data.

Figure B.1 shows a **Navgrid** example with a red line visualizing the current agent output. The small white square near the end of the red lines is the target. Based on previous outputs from the agent, the related heatmap shows the likely location of the target (higher weights are more bright and than the lower ones). Black areas in the heatmap are missing values caused by how the heatmap is rendered.

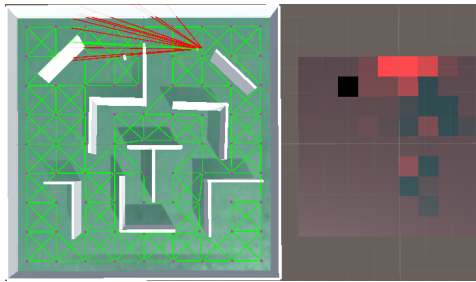


Figure B.1. An example Navgrid (left) and its weight visualized as a heatmap (right).

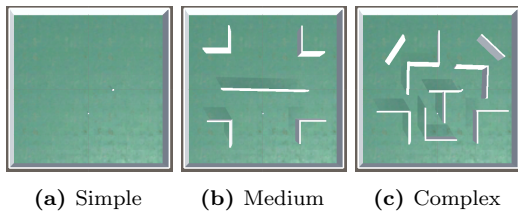


Figure B.2. The audio environments of our experiment.

695 B Audio environments

696 In the experiment we used three separate environ-
 697 ments. They can be seen in Figure B.2. Surface
 698 properties are set as Steam Audio Geometry [21],
 699 with the same absorption, scattering and transmis-
 700 sion values. Absorption affects how the low, medium
 701 and high frequencies are separately absorbed. Scat-
 702 tering affects the direction of reflection when sound
 703 is reflected from a surface. Transmission affects how
 704 much low, medium and high frequencies can pass
 705 through the surface without reflections.