

Figure 6: **Left:** An example OpenXLand environment. In this environment, the agent receives reward when the orange sphere makes contact with the blue pyramid. We see that the orange sphere is elevated, and therefore the agent must find it and use the ramps to access it. As to the blue pyramid, we do not see it because it is not there: the agent must first get the orange sphere near the black rounded cube first to spawn one. This environment also contains a grey pyramid that serves as a distraction. Importantly, if the agent brings the grey pyramid near the black rounded cube, both will disappear, making it impossible for the agent to spawn a blue pyramid and subsequently obtain its reward. **Right:** The agent’s perspective at initialization. The agent does not see the orange sphere, and so will have to search for it. However, the agent does see the grey pyramid, which could act as a distraction. As pointed out in [21], solving tasks in OpenXLand requires many skill such as navigation, exploration, experimentation, and avoiding irreversible transitions. However, since the innumerable tasks are composed of a small number of components, the agent can learn to reuse previously acquired knowledge to improve performance.

420 OpenXLand is an open-source implementation of the environment frameworks known as XLand [22]
 421 and XLand 2.0 [21]. OpenXLand generates 3D worlds from building blocks known as *metatiles*. In
 422 the world, agents are expected to solve tasks defined by *production rules*, which trigger based on the
 423 state of the agent and objects in the environment. Figure 6 shows an example of such a world and
 424 task.

425 Below, we go over specific details of the metatiles and the production rules.

426 A.1 Metatiles

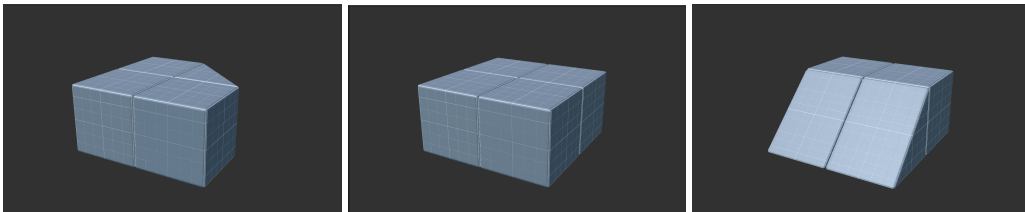


Figure 7: Three examples of metatiles. On the left, we see a platform with a rounded corner. The middle contains a square platform. The right contains a ramp, which the agent can use to traverse multi-level floor topologies.

427 The OpenXLand framework uses procedural generation to create diverse environments from sim-
 428 ple building blocks we call *metatiles*. During generation, the environment is filled with metatiles.
 429 These metatiles are designed to create interesting multi-level floor topologies: agents can use
 430 ramps to access platforms (see Figure 7). Metatiles for OpenXLand can be found in the Re-
 431 sources/GameObjects/XLand/MetaTiles subfolder of the Unity asset folder. There are 13 total

Name	Description	Sampling Weight
RampInset	2 tile wide ramp with tiles on both sides	1500
Ramp1Edge	2 tile wide ramp with tiles on one side	750
RampNoEdge	2 tile wide ramp with no tiles the sides	750
Corner	4 tile platform with one corner rounded	1000
PlatformNoEdge	4 tile platform with no internal edges	100
Platform1Edge	4 tile platform with a single internal edge	100
Platform2Edge	4 tile platform with a two internal edges at the corner	100
Platform4Edge	4 tile platform with a four internal edges	100
AirPocket	4 empty tiles	100
SmallRamp	A minimal 1 tile wide ramp	450
Air	A empty tile	45
WhiteTile	A single tile	4
AirHelper	A tile that allows an empty tile to connect to internal edges	0.0004

Table 1: Metatiles used in OpenXLand. The sampling weight is used to calculate how likely a tile is to be sampled. U3 makes it possible to easily adjust these weights to enrich environment generation dynamics.

432 metatiles distributed across two pools: a medium pool with a weight of 10, and a small pool with
433 a weight of 1. This means that the procedural generation algorithm will attempt to place medium
434 metatiles much more frequently than small metatiles. A full list of metatiles and their relative weights
435 can be found in Table 1 However, the small metatiles still serve an important role as they can fill in
436 gaps where medium metatiles no longer fit. Note that AirHelper is meant to be used only when no
437 other tiles are available.

438 **A.1.1 The Face Matching Matrix**

439 The logic for how metatiles fit together is defined by the face matching matrix. This adjacency matrix
440 allows users to define which faces can be placed next to other faces, and with what weight. U3 will
441 attempt to place metatiles at frequencies that respect both the face matching matrix frequencies and
442 the metatile pool weights, however, due to the complexity of the generation process there is a fair
443 amount of variance in the actual number of each metatile and face-face pairs placed.

444 OpenXLand uses 7 face types. Bedrock denotes the bottom of tiles, so Air and ramps cannot be
445 placed below it. Similarly, Surface denotes the top of a tile and can match with any type of face, but
446 attempts to avoid Air (0.05 weight). Air and Ground are generic face types of empty and solid faces
447 respectively. InsideEdge attempts to connect to other InsideEdge faces, to create larger platforms.
448 RampTop and RampBottom attempt to connect in a mutually exclusive manner to make sure that
449 multiple ramps can chain together in the correct way.

450 **A.1.2 Wave Function Collapse Algorithm**

451 U3 implements a 3D wave function collapse algorithm with a several optimizations. The original
452 wave function collapse algorithm tracks all possible configurations of each world voxel to calculate
453 the entropy of the voxels. We simplifi

454 the entropy of all world voxels simultaneously. In our modified version we only consider the number
455 of metatile configurations (entropy)

456 **A.2 Production Rules**

457 A key element of XLand 2.0 is *production rules*, which define compositional tasks that agents
458 are meant to complete. Agents 1) manipulate production rule objects 2) to satisfy production rule
459 conditions, which 3) trigger production rule actions that transition the state. The goal of the agent is to
460 trigger these production rules until the environment reaches a goal state and the agent receives reward.

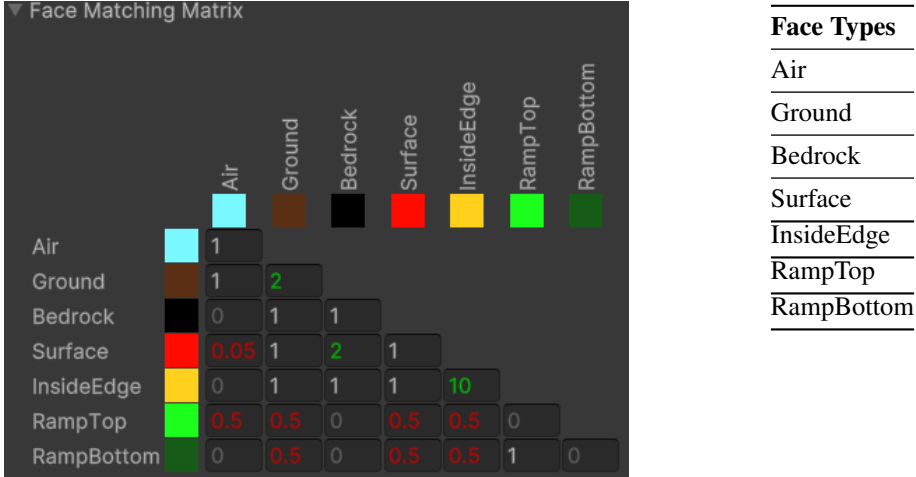


Figure 8: **Left:** The face matching matrix between face types in OpenXLand. **Right:** The face types in the adjacency matrix. Each metatile is made of tiles, which have faces, and the adjacency matrix determines which faces can be neighbors. The adjacency matrix constrains the set of permissible metatiles that can be placed at a given location during the environment generation process.

461 In Table 2, we outline the production rule objects by shape and color, as well as the production
 462 conditions and the production rule actions. Naturally, each set can be extended.

Object Shapes	Object Colors	Conditions	Actions
Cube	red	NEAR	SPAWN
Sphere	green	CONTACT	SWAP
Pyramid	blue	HOLD	REMOVE
Cylinder	purple	PICKUP	REWARD
Cone	yellow	THROW	
Rounded Box	grey	DROP	
Rounded Rectangle	white		
	etc.		

Table 2: A table containing the main components of OpenXLand’s production rule system. Our set of objects, which is specified by shape and color, is similar in nature to that of XLand 2.0 [21], though with 7 objects and over 20 colors. As for the conditions We extend those in Team et al. [21] with the new conditions “drop”, “pickup” and “throw”. Also, whereas in XLand 2.0, objects that are part of production rule conditions disappear when the condition is fulfilled, we include the possibility of objects remaining in the environment. The “spawn” action creates a new object, the “swap” action swaps the production rule condition’s subject, object or both for a new object, the “remove” action simply removes an item from the environment, and the “reward” action gives reward to the agent. Therefore, our production rules go beyond forcibly removing the condition’s objects or swapping them for new ones.

463 **A.2.1 Production Rule Objects**

464 Our production rule objects are quite similar to those of Team et al. [21]: as us, they use cubes,
 465 spheres and pyramids (not counting walls), but we also include other objects such as cylinders, cones,
 466 and rounded boxes and rectangles. Team et al. [21] uses three colors: black, purple and yellow. We
 467 see no reason a priori why there shouldn’t be more colors, so we include over 20 colors, a few of

468 which can be found in 2. This raises the unique number of object-shape pairs from 9 to over 150.
469 Naturally, we can restrict or extend the set of colors. Furthermore, although OpenXLand’s objects are
470 identified with simple colors and shapes, it is possible to use U3 to identify objects with different
471 characteristics.

472 **A.2.2 Production Rule Conditions**

473 We implement the following conditions from Team et al. [21]: “near”, “contact” and “hold“. We plan
474 on implementing the final condition, “see” in future work. On top of this, we implement three new
475 conditions: “pickup”, “throw” and “drop”, which we implement as callback functions that trigger
476 based on the agent’s behaviour. The agent can only hold (and therefore pickup, throw and drop) one
477 object at a time.

478 **A.2.3 Production Rule Actions**

479 In Team et al. [21], production rules and goal states are represented as separate concepts. We find
480 it more amenable to represent them as instances of the same concept by turning goal states into
481 production rules that reward the agent as an action. This allows us to represent production rules
482 and goal states as part of the same language. We also include in our production rule actions the
483 possibility of preserving the objects present in the production rule conditions, which can allow for
484 distinct production rules that repeatedly depend on the same object.

485 **B Implementing MemoryMaze: A Walkthrough**

486 MemoryMaze [18] differs from OpenXLand in a few keys ways: First, the world is flat, so we only
487 need a height of 1. Second, production rules are not procedurally generated, but are instead attached
488 to "goal" objects spawned throughout the world. Third, only 1 "goal" is active at a time. Fourth, there
489 is a UI element that displays the current "goal".

490 To illustrate the ability of U3 to render more naturalistic worlds, we use a low polygon forest pack
491 from the Unity asset store. Thus, we modify the original MemoryMaze by replacing the walls with
492 rocks, and the goals with forest mushrooms of various colors.

493 Our MemoryMaze starts by procedurally generating the flat world using the same metatile system as
494 OpenXLand. We want goals to be spread out and only spawn in small clearings. Thus, we define a
495 new class (ProductionRuleSpawn) that defines locations where goals can spawn, and add it to the
496 Clearing metatile itself. By adding a script into the metatile we are able to use the existing procedural
497 generation system to add structured randomness to this new spawning logic. We also include custom
498 spawn logic in the script to spawn only up to 5 goals, and make sure that each goal is a minimum
499 distance from the previously spawned goals.

500 Further, we add logic to ProductionRuleSpawn such that a new production rule is created for each
501 spawned goal. These rules all use a new NEAR condition that triggers when the agent is near the
502 production rule’s subject (in this case, the goal object). Because of the modular design of U3, this
503 new condition can now be used in any newly generated OpenXLand environments as well. We also
504 create a new REWARD_TOGGLE action that is a compound action of a REWARD plus new logic
505 that deactivates the current production rule, and activates a different production rule at random. To
506 implement this logic we also add a new flag to production rules that allows them to be turned on and
507 off at run-time. In this way, as new environments are added to the U3 ecosystem, we expect that new
508 and more powerful features will also be added to the base OpenXLand environments.

509 The final missing piece of the MemoryMaze environment is the UI element that shows the current
510 goal. We adopt the same UI as the original with a colored border around the image. Unity is a
511 powerful 3D game engine with native support for 3D rendering, UI, sprites, materials, etc, so creating
512 this UI took only 1 function call into Unity’s API (GUI.DrawTexture).



Figure 9: MemoryMaze environment observation. Our version of MemoryMaze uses naturalistic assets. We hope that harnessing the variety of assets available will help to increase agents’ overall adaptability.

513 C Dataset Details

514 In this section, we go over the details of the datasets. We split the datasets into 6 world datasets 3,
515 and 6 production rule datasets 4. We design the datasets with overlapping parameters to facilitate
516 smooth curriculum learning. We filter the world datasets to include only those worlds that have a
517 largest accessible area of over 50% of the total number of tiles in the X-Z plane.

518 We find the largest accessible area by first constructing a local, directed connectivity graph using
519 the specific configuration of each voxel. For simplicity, we consider only the cardinal connections
520 between voxels, and only connections in the width-length plane (we ignore connections within a
521 single column of voxels). Two voxels are connected bi-directionally if they are at the same terrain
522 height and higher voxels are connected unidirectionally to lower tiles - representing the agent’s ability
523 to fall to the lower level. Ramps serve as the sole method for connections from a lower layer to the
524 next highest layer. Once the connectivity graph is constructed, the largest accessible area is found by
525 computing the largest strongly connected sub-component.

526 Distractors’ positions in the production rule chains heavily influence their impact. For example, if the
527 agent has already completed the step that requires a Yellow Cube, then removing the Yellow Cube is
528 merely a distraction. However, if the Yellow Cube is still required, then its removal creates a dead
529 end. Due to this complexity, we do not directly consider distractors vs dead ends in our difficulty
530 measure, but note that as more distractors are added the probability of dead ends also increases.

531 D Training Results

532 To demonstrate the learnability of our environment, we randomly select a task from the easy dataset
533 and applied the Soft Actor-Critic (SAC) [11] method. Image observations are processed using a
534 convolutional neural network for both the actor and Q networks, each comprising approximately 2
535 million parameters. The SAC agent is trained for 75 million environment steps, utilizing a replay
536 buffer that stored 1 million transitions. While these parameter choices may appear modest compared
537 to those in the Ada paper [21] and may not fully exploit emergent meta-learning capabilities, our
538 primary aim is to underscore the learnability of our environment, not to recreate their agent.

Name	Width and Length		Height	
	\bar{x}	σ	\bar{x}	σ
easy_low	5	2	1	1
easy_high	5	2	3	2
medium_low	10	3	1	1
medium_high	10	3	3	2
hard_low	15	5	1	1
hard_high	15	5	3	2

Table 3: World Datasets. These values were used to define Gaussian distributions used to sample parameters during dataset construction. We note that the harder we consider environments to be, the bigger they are.

Name	Chain Length and Initial Objects		Distractors	
	\bar{x}	σ	\bar{x}	σ
short_few	1	1	0	1
short_many	1	1	1	0.5
middle_few	3	1.5	0	1
middle_many	3	1.5	2	1.5
long_few	5	3	0	1
long_many	5	3	3	2.5

Table 4: Production Rule Datasets. These values were used to define the Gaussian distribution used to sample parameters during dataset construction. Since we sample the production rule chain lengths from distributions, the datasets of shorter chains will still contain long chains, although the dataset will be dominated by shorter, less complicated tasks. As the average length of chains increases throughout the datasets, the number of distractors also increases.

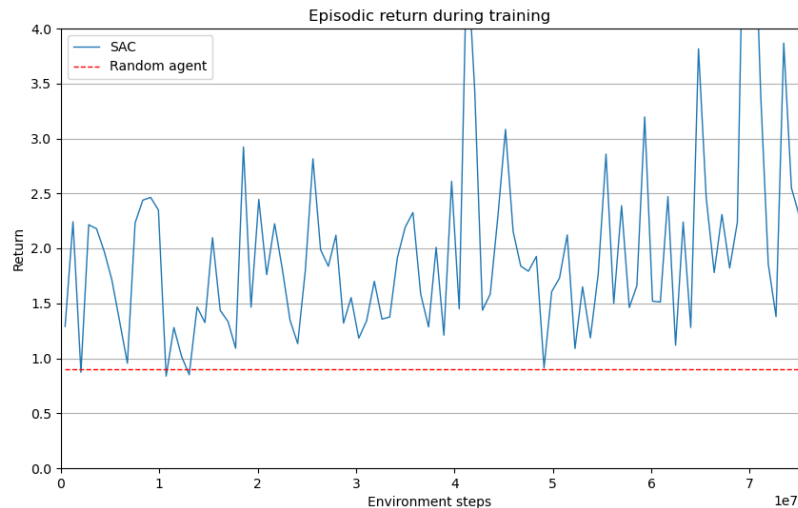


Figure 10: Episodic return during single environment training of the SAC agent compared to a random agent. We train the SAC agent for 75 million environment steps. This required approximately 48 hours on a single A100 GPU, a small fraction of the 53760 TPUv3 device hours it took to train Ada [21]. The episodic return for the SAC agent shows significant fluctuations but generally demonstrates a higher return compared to the random agent baseline, represented by the red dashed line.

539 Figure 10 illustrates that the episodic return of our learning agent increases throughout training,
 540 surpassing a random agent baseline’s performance. We also extend our evaluation to the entire easy
 541 dataset, conducting 8 trials per episode on different tasks. Interestingly, our learning agent exhibits
 542 minimal improvement over random agents in this broader context, indicating the limitations of
 543 conventional methods. We observe similar outcomes when applying the method to medium datasets,
 544 where random agents fail to accrue any rewards, underscoring the challenge of reward collection
 545 within these datasets.

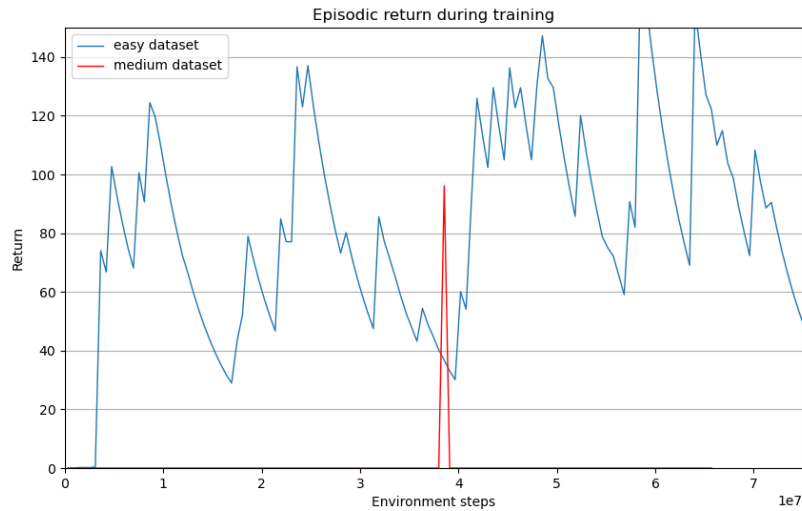


Figure 11: Episodic return during training using SAC on the whole dataset with 8 trials per episode and changing the task each episode. Collecting reward on the medium dataset is hard showing the limitations of the conventional methods. Collecting reward on easy dataset is easier, but meta-learning requires more advanced methodologies.

546 E Resources

547 Most of the resources for this project were spent on implementing the baselines. The CPU compute
 548 spent to generate the envs amounts to 200 hours. The GPU compute spent to train the baselines
 549 amounts to around 140 hours of NVIDIA A100-SXM4-80GB.

550 F URL to Dataset and Metadata

551 **Datasets:** https://huggingface.co/datasets/cerc-aai/u3_datasets

552 **U3 source code:** <https://github.com/CERC-AAI/u3/>

553 **Croissant Metadata:** [https://huggingface.co/api/datasets/cerc-aai/u3_datasets/](https://huggingface.co/api/datasets/cerc-aai/u3_datasets/croissant)
 554 croissant

555 We plan to fully open-source the code for U3, including the Unity environment framework, the
 556 python training code and the pipeline. Using U3 still requires installing Unity, which impinges on
 557 reproducibility. As stated in Section 7, we plan on containerizing the environment to democratize
 558 access to U3.

559 **G Datasheet**

560 **G.1 Motivation**

- 561 1. **For what purpose was the dataset created?** *U3 and its associated OpenXLand datasets*
562 *were created to promote the open development of foundation models for RL and meta-RL.*
- 563 2. **Who created the dataset and on behalf of which entity?** *U3 and its datasets were created*
564 *by the researchers listed in the author list.*
- 565 3. **Who funded the creation of the dataset?** *The main funding bodies include the Canada*
566 *Excellence Research Chairs Program, as well as the resources of the INCITE program*
567 *award “Scalable Foundation Models for Transferable Generalist AI” provided by Oak*
568 *Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is*
569 *supported by the Office of Science of the U.S. Department of Energy under Contract No.*
570 *DE-AC05-00OR22725*

571 **G.2 Distribution**

- 572 1. **Will the dataset be distributed to third parties outside of the entity (e.g., company,**
573 **institution, organization) on behalf of which the dataset was created?** *The U3 framework*
574 *and its datasets will be publicly available.*
- 575 2. **How will the dataset be distributed (e.g., tarball on website, API, GitHub)?** *The datasets*
576 *will be distributed through Hugging Face, while the code will be made available on Github.*
- 577 3. **Have any third parties imposed IP-based or other restrictions on the data associated**
578 **with the instances?** *No.*
- 579 4. **Do any export controls or other regulatory restrictions apply to the dataset or to**
580 **individual instances?** *No.*

581 **G.3 Maintenance**

- 582 1. **Who will be supporting/hosting/maintaining the dataset?** *The CERC-AAI Lab at Univer-*
583 *sité de Montréal will maintain the U3 framework and its datasets.*
- 584 2. **How can the owner/curator/manager of the dataset be contacted (e.g., email ad-**
585 **dress)?** *The managers of the dataset can be contacted at connor.brennan@mila.quebec,*
586 *andrew.williams@mila.quebec and irina.rish@gmail.com*
- 587 3. **Is there an erratum?** *No. We will release one in the future if necessary.*
- 588 4. **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**
589 **instances)?** *Yes, the datasets will be updated if necessary.*
- 590 5. **If the dataset relates to people, are there applicable limits on the retention of the data**
591 **associated with the instances (e.g., were the individuals in question told that their data**
592 **would be retained for a fixed period of time and then deleted?)** *N/A*
- 593 6. **Will older versions of the dataset continue to be supported/hosted/maintained?** *N/A*
- 594 7. **If others want to extend/augment/build on/contribute to the dataset, is there a mecha-**
595 **nism for them to do so?** *U3 welcomes contributors on GitHub. Should the need arise to*
596 *extend the datasets on Hugging Face, we will establish a procedure for doing so.*

597 **G.4 Composition**

- 598 1. **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
599 **countries?)** *The instances that comprise the datasets are JSON datafiles that describe the*
600 *necessary information to instantiate RL environments in Unity.*
- 601 2. **How many instances are there in total (of each type, if appropriate)?** *Six files of one*
602 *million environment instances, six files of one million production rule instances, for a total*
603 *of twelve million instances spread across 12 files.*

- 604 3. **Does the dataset contain all possible instances or is it a sample of instances from a**
605 **larger set?** *Is it a subset of the effectively infinite number of instances that U3 can generate.*
- 606 4. **Is there a label or target associated with each instance?** *Each instance contains the label*
607 *information necessary to instantiate a Unity environment.*
- 608 5. **Is any information missing from individual instances?** *No.*
- 609 6. **Are there recommended data splits (e.g., training, development/validation, testing)?**
610 *No.*
- 611 7. **Are there any errors, sources of noise, or redundancies in the dataset?** *The instances*
612 *are sampled with replacement from a procedural generation process, so yes there can be*
613 *redundancies.*
- 614 8. **Is the dataset self-contained, or does it link to or otherwise rely on external resources**
615 **(e.g., websites, tweets, other datasets)?** *The datasets rely on access to a Unity instance*
- 616 9. **Does the dataset contain data that might be considered confidential?** *No.*
- 617 10. **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**
618 **threatening, or might otherwise cause anxiety?** *No.*

619 G.5 Collection Process

- 620 1. **How was the data associated with each instance acquired?** *We procedurally generated*
621 *environments in Unity and collected the necessary information to recreate them using the*
622 *U3 framework.*
- 623 2. **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus**
624 **or sensor, manual human curation, software program, software API)?** *We used the*
625 *Unity Engine and Python, as well as GPUs for training RL baselines.*
- 626 3. **Who was involved in the data collection process (e.g., students, crowdworkers, con-**
627 **tractors) and how were they compensated (e.g., how much were crowdworkers paid)?**
628 *Students, paid through the funding sources described above.*
- 629 4. **Does the dataset relate to people?** *No.*
- 630 5. **Did you collect the data from the individuals in question directly, or obtain it via third**
631 **parties or other sources (e.g., websites)?** *N/A*

632 G.6 Uses

- 633 1. **Has the dataset been used for any tasks already?** *No, the data has not been used, other*
634 *than for running baselines.*
- 635 2. **What (other) tasks could the dataset be used for?** *The datasets could be used for*
636 *RL-related tasks as standalone or as part of a blended dataset*
- 637 3. **Is there anything about the composition of the dataset or the way it was collected and**
638 **preprocessed/cleaned/labeled that might impact future uses?** *The use of U3 and its*
639 *datasets relies on access to the Unity Engine.*
- 640 4. **Are there tasks for which the dataset should not be used?** *No.*

641 H Author statement

642 We bear all responsibility in case of violation of rights.

643 We use the MIT License for our data.