

## 476 A What Is Inside the Datasets?

477 Every dataset is repacked into HDF5 files similar to Fu et al. (2020). The data keys are described  
478 in Table 3; along to the typical  $(s_t, a_t, r_t, d_t)$  tuples, the metadata is also provided as the datasets’  
479 attributes with a comprehensive information about specific trajectories similar to Hambro et al.  
480 (2022b). The re-packing script is provided at [https://github.com/tinkoff-ai/katakomba/  
481 scripts/generate\\_small\\_dataset.py](https://github.com/tinkoff-ai/katakomba/scripts/generate_small_dataset.py).

Table 3: The re-packed datasets constitute of transformed data from Hambro et al. (2022b). Dissimilar  
the the large scale dataset, the repacked data is now in the format familiar to the ORL practitioners.  
We also save the metadata for each trajectory, for a comprehensive description, please, see Appendix  
F in Hambro et al. (2022b).

Name	Type	Shape	Description
tty_chars	np.uint8	[B, T, H, W]	$s_t$ : The on-screen characters (default screen size 80 x 24).
tty_colors	np.int8	[B, T, H, W]	$s_t$ : The on-screen colors for each character.
tty_cursor	np.int16	[B, T, 2]	$s_t$ : The coordinates of the on-screen cursor.
actions	np.uint8	[B, T]	$a_t$ : The NLE actions the player made in response to the $s_t$ .
rewards	np.int32	[B, T]	$r_t$ : The difference between in-game scores at states $s_t$ and $s_{t-1}$ . Note that this was used in all implementations of the algorithms provided in Hambro et al. (2022b). We also found that without this reward-shaping, all offline RL algorithms failed completely.
done	np.uint8	[B, T]	$d_t$ : An indicator whether the current state is the last one in the trajectory.

## 482 B License

483 Our codebase and repacked datasets are released under the NETHACK GENERAL PUBLIC LI-  
484 CENSE. The original NetHack Learning environment (Küttler et al., 2020) and large-scale datasets  
485 (Hambro et al., 2022b) are also released under NETHACK GENERAL PUBLIC LICENSE.

## 486 C General Ethic Conduct and Potential Negative Societal Impact

487 To the best of our knowledge, our work does not present any direct potential negative societal impact.  
488 As of the general ethic conduct, we believe that the most relevant issue to be discussed is the  
489 "Consent to use or share the data". Our work is largely built upon both the NetHack Learning  
490 Environment (Küttler et al., 2020) and the corresponding large-scale dataset (Hambro et al., 2022b),  
491 and as already described in the Appendix B both are distributed under the NETHACK GENERAL  
492 PUBLIC LICENSE that explicitly allows for re-usage and re-distribution.

493 **D Resources and Statistics**

494 We used 64 separated computational nodes with 1xA100, 14CPU, 128GB RAM, and the NVMe as  
 495 long-term storage for all our experiments. All the values reported in the paper were also obtained  
 496 under this configuration. One can also find more detailed information inside the Weights&Biases  
 497 logs in the code repository.

Table 4: Scores used for Normalization. You can also find them at <https://github.com/tinkoff-ai/katakomba/katakomba/utis/scores.py>. For other statistics, please, see Table 2 in the main text.

Tasks	Minimum Score	Maximum Score	Mean Score
<b>Base (Role-Centric)</b>	-	-	-
<u>arc</u> -hum-neu	0.0	138103.0	6636.44
<u>bar</u> -hum-neu	0.0	292342.0	17836.68
<u>cav</u> -hum-neu	0.0	258978.0	12113.87
<u>hea</u> -hum-neu	0.0	64337.0	4068.27
<u>kni</u> -hum-law	0.0	419154.0	14137.06
<u>mon</u> -hum-neu	0.0	171224.0	17456.05
<u>pri</u> -hum-neu	0.0	114269.0	7732.69
<u>ran</u> -hum-neu	0.0	54874.0	8067.99
<u>rog</u> -hum-cha	0.0	68628.0	4818.20
<u>sam</u> -hum-law	0.0	155163.0	11009.36
<u>tou</u> -hum-neu	0.0	59484.0	4211.47
<u>val</u> -hum-neu	16.0	313858.0	18624.77
<u>wiz</u> -hum-neu	0.0	71709.0	5323.48
<b>Extended (Race-Centric)</b>	-	-	-
<u>pri</u> -elf-cha	0.0	83744.0	7109.35
<u>ran</u> -elf-cha	0.0	66690.0	9014.18
<u>wiz</u> -elf-cha	0.0	71664.0	5005.16
<u>arc</u> -dwa-law	0.0	83496.00	5445.69
<u>cav</u> -dwa-law	0.0	161682.0	11893.48
<u>val</u> -dwa-law	0.0	1136591.0	23473.61
<u>arc</u> -gno-neu	0.0	110054.0	5316.57
<u>cav</u> -gno-neu	0.0	142460.0	10083.06
<u>hea</u> -gno-neu	0.0	69566.0	3783.93
<u>ran</u> -gno-neu	0.0	58137.0	6965.04
<u>wiz</u> -gno-neu	0.0	37376.0	4317.51
<u>bar</u> -orc-cha	0.0	164296.0	17594.38
<u>ran</u> -orc-cha	3.0	69244.0	7608.48
<u>rog</u> -orc-cha	0.0	54892.0	4897.69
<u>wiz</u> -orc-cha	0.0	40871.0	5016.74
<b>Complete (Alignment-Centric)</b>	-	-	-
<u>arc</u> -hum-law	2.0	84823.0	5826.35
<u>cav</u> -hum-law	0.0	156966.0	12462.82
<u>mon</u> -hum-law	7.0	190783.0	16091.57
<u>pri</u> -hum-law	0.0	99250.0	6847.99
<u>val</u> -hum-law	0.0	428274.0	26103.03
<u>bar</u> -hum-cha	0.0	164446.0	18228.11
<u>mon</u> -hum-cha	0.0	223997.0	18353.30
<u>pri</u> -hum-cha	0.0	58367.0	8262.56
<u>ran</u> -hum-cha	3.0	62599.0	8378.50
<u>wiz</u> -hum-cha	0.0	55185.0	5316.82

498 **E Hyperparameters**

499 For all algorithms, hyperparameters have been reused from previous work whenever possible. For  
 500 BC, CQL, and IQL reference values, see Appendix I.4 in the Hambro et al. (2022b). For AWAC,  
 501 hyperparameters from IQL were reused due to the very similar policy updating scheme. For REM,  
 502 hyperparameters were taken from the original work (see Agarwal et al. (2020)).

503 As in Hambro et al. (2022b), and in contrast to the original CQL implementation, we multiply the TD  
 504 loss by the  $\alpha$  coefficient instead of the CQL loss, as we observed better results with such a scheme.  
 505 We performed a search for  $\alpha \in [0.0001, 0.0005, 0.001, 0.05, 0.01, 0.05, 0.1, 0.5, 1.0]$  with best value  
 506  $\alpha = 0.0001$ .

Table 5: BC hyperparameters.

Parameter	Value
optimizer	AdamW (Kingma & Ba, 2014; Loshchilov & Hutter, 2017)
training iterations	500000
batch size	64
sequence length	16
learning rate	3e-4
weight decay	0.0
state encoder	Chaotic-Dwarven-GPT-5(Hambro et al., 2022a,b)
LSTM hidden dim	2048
LSTM layers	2
LSTM dropout	0.0
use previous action	True

Table 6: CQL hyperparameters. Note that in our implementation, the  $\alpha$  coefficient multiplies the TD loss.

Parameter	Value
optimizer	AdamW (Kingma & Ba, 2014; Loshchilov & Hutter, 2017)
training iterations	500000
batch size	64
sequence length	16
learning rate	3e-4
weight decay	0.0
state encoder	Chaotic-Dwarven-GPT-5(Hambro et al., 2022a,b)
LSTM hidden dim	2048
LSTM layers	2
LSTM dropout	0.0
use previous action	True
tau ( $\tau$ )	5e-3
gamma ( $\gamma$ )	0.999
reward clip range	[-10.0, 10.0]
alpha ( $\alpha$ )	1e-4

Table 7: IQL hyperparameters.

Parameter	Value
optimizer	AdamW (Kingma & Ba, 2014; Loshchilov & Hutter, 2017)
training iterations	500000
batch size	64
sequence length	16
learning rate	3e-4
weight decay	0.0
state encoder	Chaotic-Dwarven-GPT-5(Hambro et al., 2022a,b)
LSTM hidden dim	2048
LSTM layers	2
LSTM dropout	0.0
use previous action	True
tau ( $\tau$ )	5e-3
gamma ( $\gamma$ )	0.999
reward clip range	[-10.0, 10.0]
expectile	0.8
temperature	1.0
advantage clip max	100

Table 8: AWAC hyperparameters.

Parameter	Value
optimizer	AdamW (Kingma & Ba, 2014; Loshchilov & Hutter, 2017)
training iterations	500000
batch size	64
sequence length	16
learning rate	3e-4
weight decay	0.0
state encoder	Chaotic-Dwarven-GPT-5(Hambro et al., 2022a,b)
LSTM hidden dim	2048
LSTM layers	2
LSTM dropout	0.0
use previous action	True
tau ( $\tau$ )	5e-3
gamma ( $\gamma$ )	0.999
reward clip range	[-10.0, 10.0]
temperature	1.0
advantage clip max	100

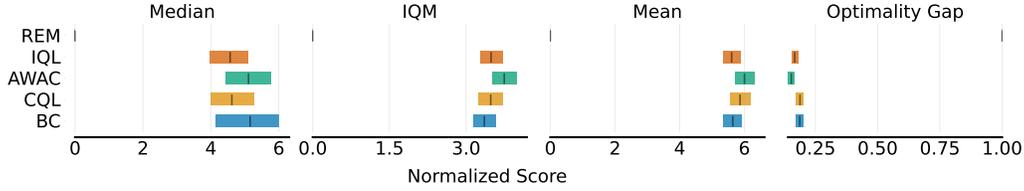
Table 9: REM hyperparameters.

Parameter	Value
optimizer	AdamW (Kingma & Ba, 2014; Loshchilov & Hutter, 2017)
training iterations	500000
batch size	64
sequence length	16
learning rate	3e-4
weight decay	0.0
state encoder	Chaotic-Dwarven-GPT-5(Hambro et al., 2022a,b)
LSTM hidden dim	2048
LSTM layers	2
LSTM dropout	0.0
use previous action	True
tau ( $\tau$ )	5e-3
gamma ( $\gamma$ )	0.999
reward clip range	[-10.0, 10.0]
ensemble heads	200.0

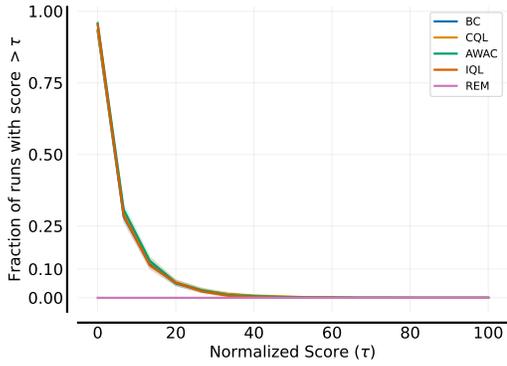
507 **F Results per Benchmark Categories**

508 In this section, we report the results stratified by the introduced categories. If one is willing to  
 509 inspect specific datasets, we organized all training logs into Weights&Biases public reports, found at  
 510 <https://wandb.ai/tlab/NetHack/reports>.

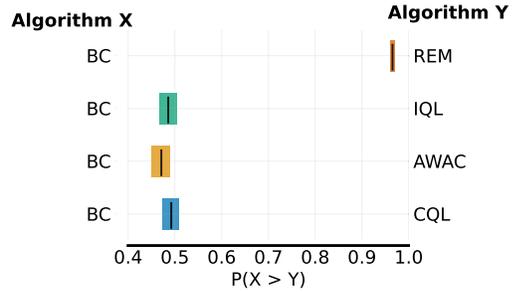
511 Note that one can find all the evaluation scores (for more than one checkpoint) within the runs and use  
 512 them for any evaluation tools of interest. Also, we provide convenient scripts for constructing RLIable  
 513 (Agarwal et al., 2021) graphs based on the provided runs that can be configured for one’s purposes as  
 514 well (see [https://github.com/tinkoff-ai/katakomba/scripts/rliable\\_report.py](https://github.com/tinkoff-ai/katakomba/scripts/rliable_report.py)).



(a) Bootstrapped point estimates.

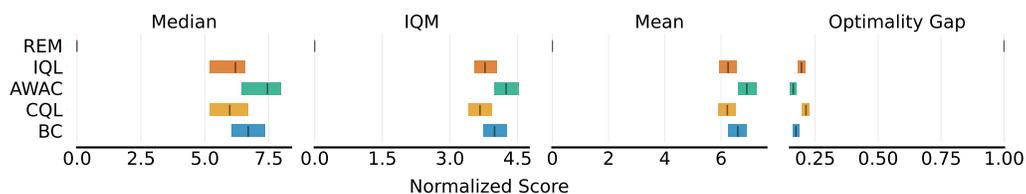


(b) Performance profiles.

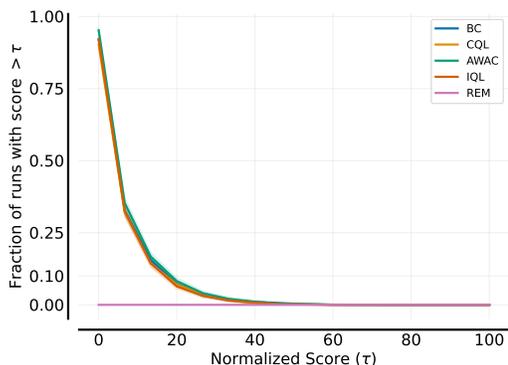


(c) Probability of improvement of BC to other algorithms.

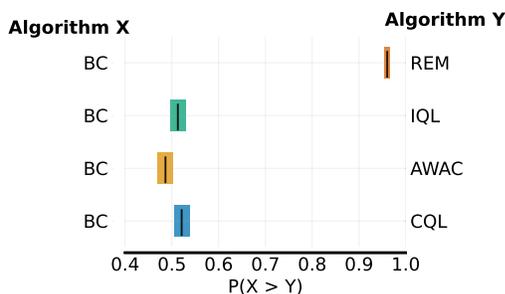
Figure 3: Normalized performance under the Katakomba benchmark for Base datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 1950 points for constructing these graphs.



(a) Bootstrapped point estimates.

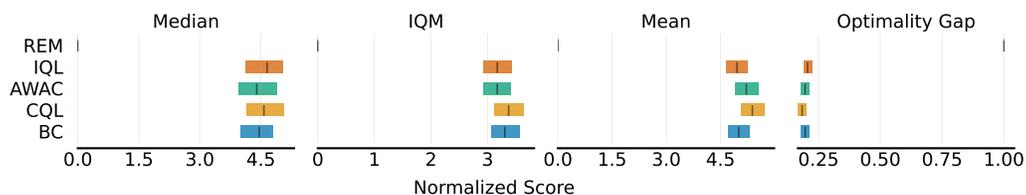


(b) Performance profiles.

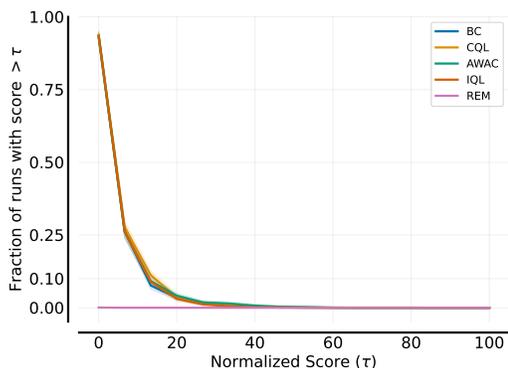


(c) Probability of improvement of BC to other algorithms.

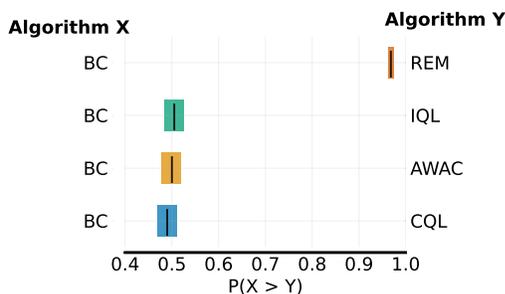
Figure 4: Normalized performance under the Katakomba benchmark for Extended datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 2250 points for constructing these graphs.



(a) Bootstrapped point estimates.

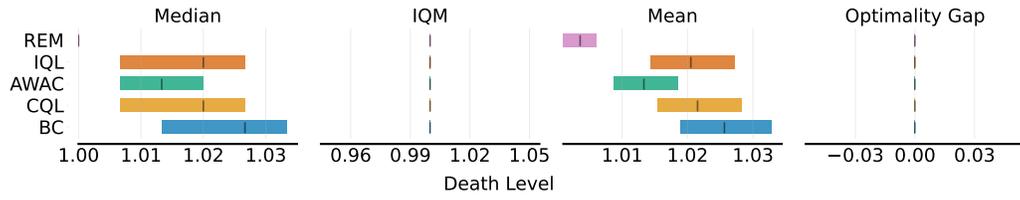


(b) Performance profiles.

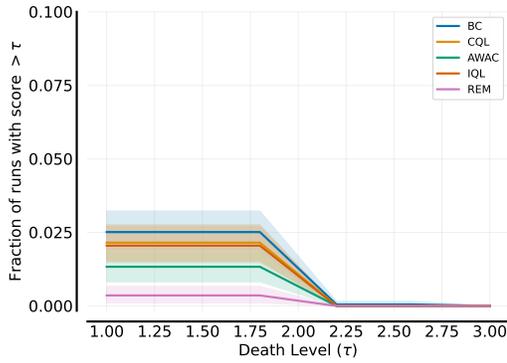


(c) Probability of improvement of BC to other algorithms.

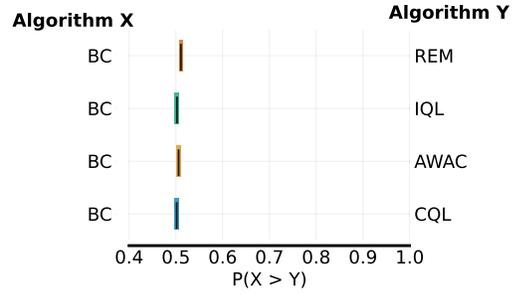
Figure 5: Normalized performance under the Katakomba benchmark for Complete datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 1500 points for constructing these graphs.



(a) Bootstrapped point estimates.

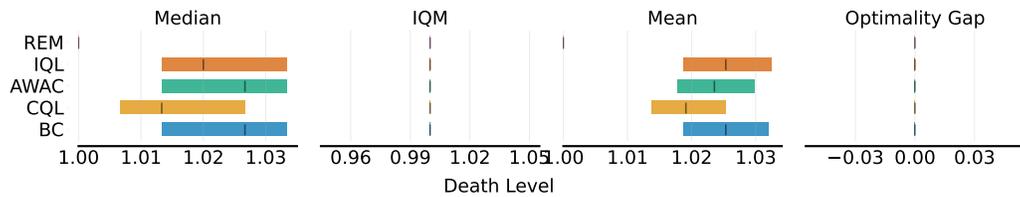


(b) Performance profiles.

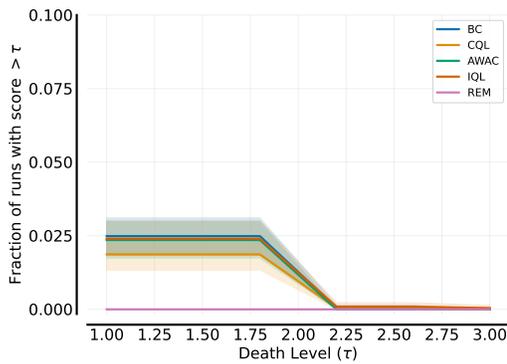


(c) Probability of improvement of BC to other algorithms.

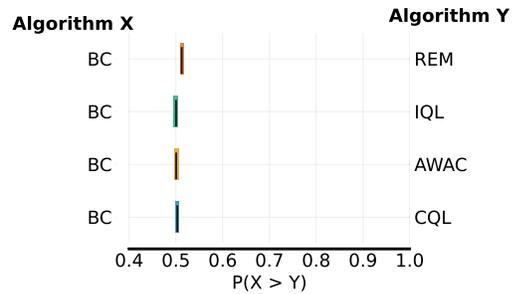
Figure 6: Death levels under the Katakomba benchmark for Base datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 1950 points for constructing these graphs.



(a) Bootstrapped point estimates.

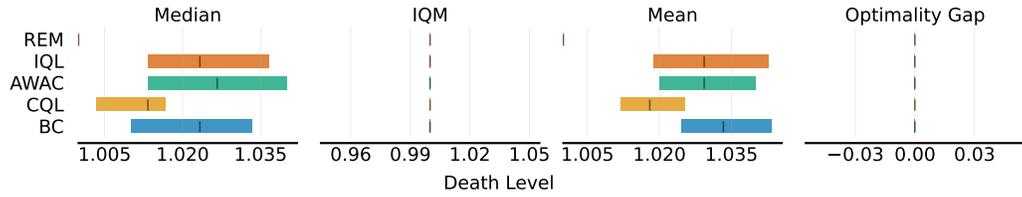


(b) Performance profiles.

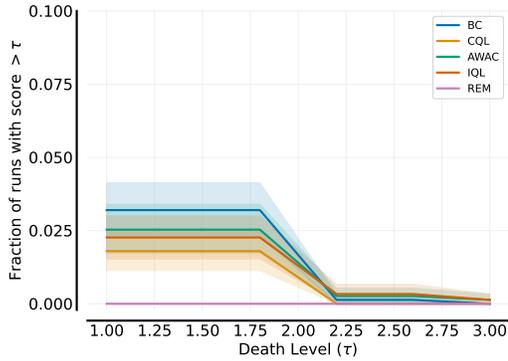


(c) Probability of improvement of BC to other algorithms.

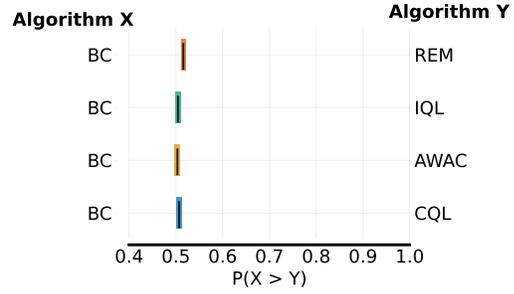
Figure 7: Death level under the Katakomba benchmark for Extended datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 2250 points for constructing these graphs.



(a) Bootstrapped point estimates.

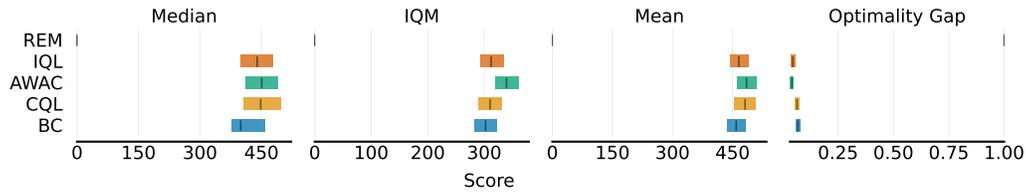


(b) Performance profiles.

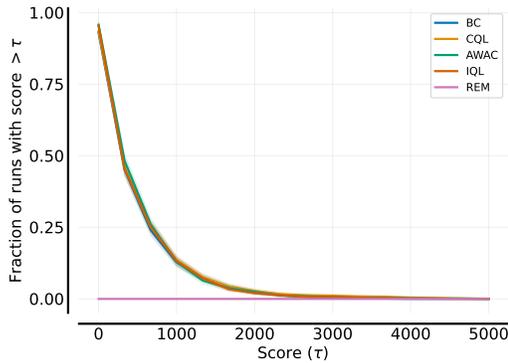


(c) Probability of improvement of BC to other algorithms.

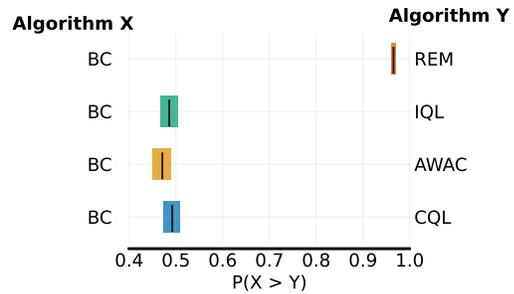
Figure 8: Death levels under the Katakomba benchmark for Complete datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 1500 points for constructing these graphs.



(a) Bootstrapped point estimates.

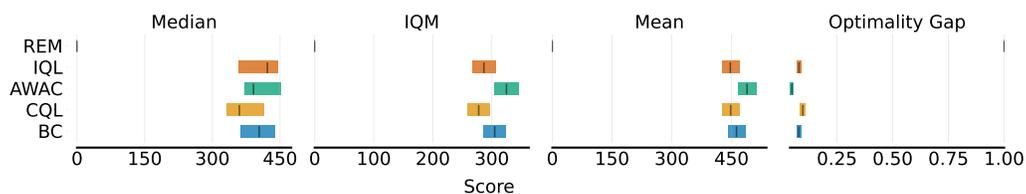


(b) Performance profiles.

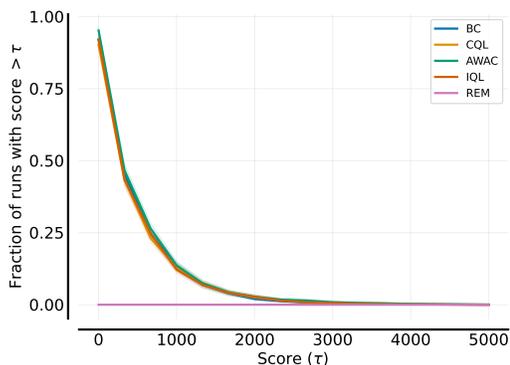


(c) Probability of improvement of BC to other algorithms.

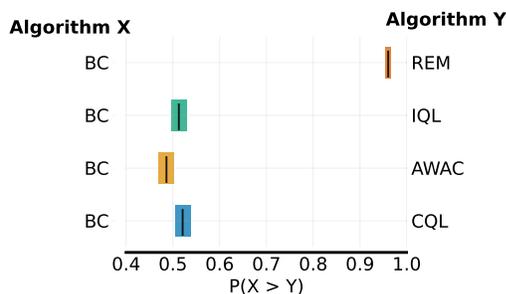
Figure 9: Unnormalized in-game score under the Katakomba benchmark for Base datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 1950 points for constructing these graphs.



(a) Bootstrapped point estimates.

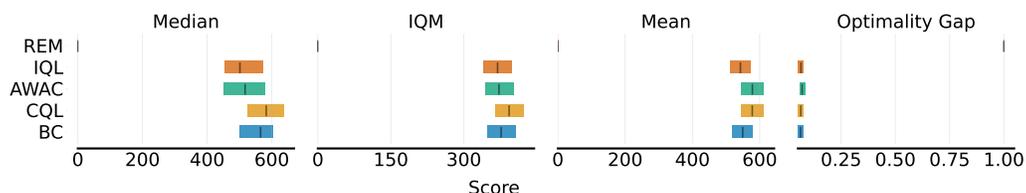


(b) Performance profiles.

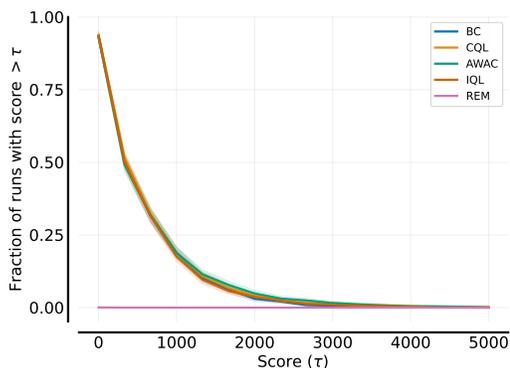


(c) Probability of improvement of BC to other algorithms.

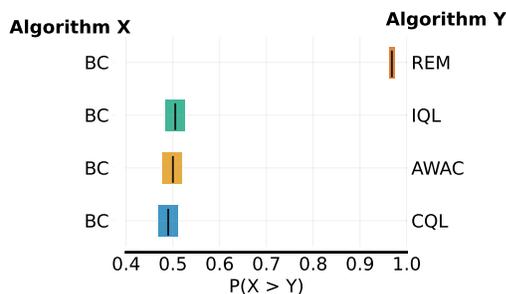
Figure 10: Unnormalized in-game score under the Katakomba benchmark for Extended datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 2250 points for constructing these graphs.



(a) Bootstrapped point estimates.



(b) Performance profiles.



(c) Probability of improvement of BC to other algorithms.

Figure 11: Unnormalized in-game score under the Katakomba benchmark for Complete datasets. Each algorithm was run for three seeds and evaluated over 50 episodes resulting in 1500 points for constructing these graphs.