
SynIL: Leveraging Synergy for Offline Imitation Learning from Imperfect Demonstration Datasets

Anonymous authors
Paper under double-blind review

Abstract

Imitation learning has undergone significant evolution with the advent of deep learning. Deep neural networks enable the learning of complex policies directly from demonstrations, without relying on traditional handcrafted feature approaches. However, deep learning-based imitation learning requires high-quality demonstrations, as the policies are directly trained on data. Thus, factors such as task difficulty and expert proficiency can lead to contamination of non-optimal demonstrations in reality, resulting in decreased performance. Previous studies have explored methods for evaluating the quality of demonstrations to identify optimal samples of demonstration. However, they often require the hand-selection of expert data in advance, which becomes increasingly challenging as datasets grow larger. This study proposes a method for automated evaluation of demonstration quality based on synergy, a low-dimensional structure observed in biological systems. The proposed method quantifies the degree of synergy manifestation as rewards to perform offline reinforcement learning. Since synergy is reported to relate to proficiency, we expect this to work as an indicator of the motion quality of demonstrations. Results demonstrated that the synergy-based rewards correlated with true rewards, and synergy-based imitation learning outperformed behavior cloning and even offline reinforcement learning with true rewards in some cases. Those results will offer a new framework for enhancing imitation learning systems with demonstrations in which not every sample is optimal. The proposed method will also contribute to computational neuroscience, as well as robotics and machine learning.

1 Introduction

Imitation learning, a method for mimicking expert demonstrations, has gained attention as a new learning paradigm that can potentially outperform manual programming requiring considerable time and a high level of expertise in coding and control. Unlike traditional manual programming methods, imitation learning is a data-driven approach, allowing learning without the need for complex design, programming skills, and specialized knowledge. Imitation learning is highly regarded as a more realistic approach in fields such as robotics, healthcare, and autonomous driving (Kim et al., 2024; Pan et al., 2017). In the field of robotics, imitation learning has successfully automated numerous actions in various robots, including Mobile ALOHA (Fu et al., 2024). Recent robotic foundation models are also based on the framework of imitation learning (Kawaharazuka et al., 2025).

In spite of its potential, the performance of imitation learning is adversely affected by low quality teaching data. A large amount of optimal demonstration data is required to ensure the stable performance of imitation learning. However, due to factors such as the difficulty of the task, the skill level of the expert, and issues like fatigue and concentration when collecting demonstration data, the data obtained may include suboptimal data. As the dataset grows larger, the likelihood of including incomplete data and noise increases. Presence of suboptimal or poor-quality demonstrations gets imitation learning to mimic these poor-quality demonstrations, resulting in decreased performance.

054 Evaluating the quality of demonstrations will overcome the above issue. For example, De-
055 moDICE (Kim et al., 2022), ISWBC (Li et al., 2024), and DWBC (Xu et al., 2022) pre-divide
056 the training data into optimal data D_E from experts and D_U with unknown proficiency,
057 enhancing imitation learning performance by focusing more on expert data or adjusting the
058 distributional distance with non-optimal data. Additionally, CLUE (Liu et al., 2023) and
059 ORIL (Zolna et al., 2020) address this issue by converting distributional distances with op-
060 timal data into rewards and conducting offline reinforcement learning, which considers these
061 rewards. This approach has the potential to enable to learn behaviors beyond the (possibly
062 poor) quality of the demonstrations. However, in all these studies, it is necessary for humans
063 to pre-determine the expert data. Furthermore, when it comes to qualitative evaluation, it
064 may be challenging to discern which instructional actions are optimal depending on the task.
065 Therefore, a method to automatically and quantitatively evaluate optimality is needed.

066 To realize automated demonstration-quality assessment, we focus on the concept of syn-
067 ergy, which is proposed in neuroscience. Synergy refers to the low-dimensional structure of
068 coordinated muscles and joints, which are often observed in human proficient movements
069 (Latash et al., 2014; Turvey, 1990). Many studies have reported that the degree of synergy
070 manifestation corresponds to the proficiency of movements (Zaal et al., 1999; Gentner et al.,
071 2010). This phenomenon has also been observed from reinforcement learning in locomotion
072 (Chai & Hayashibe, 2020) and reaching movements (Han et al., 2021). Synergies and their
073 degree of manifestation can be obtained algorithmically. Thus, we expect that synergy can
074 be used for automated demonstration-quality assessment.

075 In this study, we propose a method of imitation learning with automatic assessment of
076 demonstration quality. We propose a method to compute rewards from unlabelled demon-
077 strations based on synergy and use them in offline reinforcement learning. Our contributions
078 are as follows:

- 079 • We propose a method to algorithmically compute the quality of demonstrations
080 based on synergy and use it as reward for learning.
- 081 • This study is the first application of the concept of synergy to assess the quality of
082 demonstration movements in imitation learning.
- 083 • We successfully improved imitation learning performance across various datasets,
084 including real human demonstrations.

086 2 Methodology

087 2.1 Overview

088 This section describes the proposed method, SynIL, which derives its name from synergy-
089 based imitation learning. This method allows a policy to imitate demonstrations by per-
090 forming offline reinforcement learning based on rewards computed from the degree of synergy
091 manifestation. In SynIL, demonstrations do not have to include rewards; instead, rewards
092 are computed from demonstrations using synergy, a concept proposed in neuroscience that
093 quantifies the degree of movement coordination. Unlike most conventional methods, SynIL
094 does not require manual evaluation of demonstrations, except for setting several hyperpa-
095 rameters.

096 An overview of SynIL is provided in Figure 1. First, the degree of synergy manifestation is
097 quantified as a synergy score ξ . Next, the synergy scores ξ are converted into rewards by
098 using self-supervised reward regression (SSRR) (Chen et al., 2021). Finally, offline reinforce-
099 ment learning is performed using the synergy-based rewards associated with state-action
100 pairs.

101 2.2 Acquisition of Synergy Scores based on Synergy Manifestation

102 As the first step of SynIL, the quality of demonstration is quantified based on synergy.
103 Among various types of synergy, we employ the spatial synergy, which represents low-
104 dimensional coordination patterns in multi-dimensional spatial data. The spatial synergy is
105

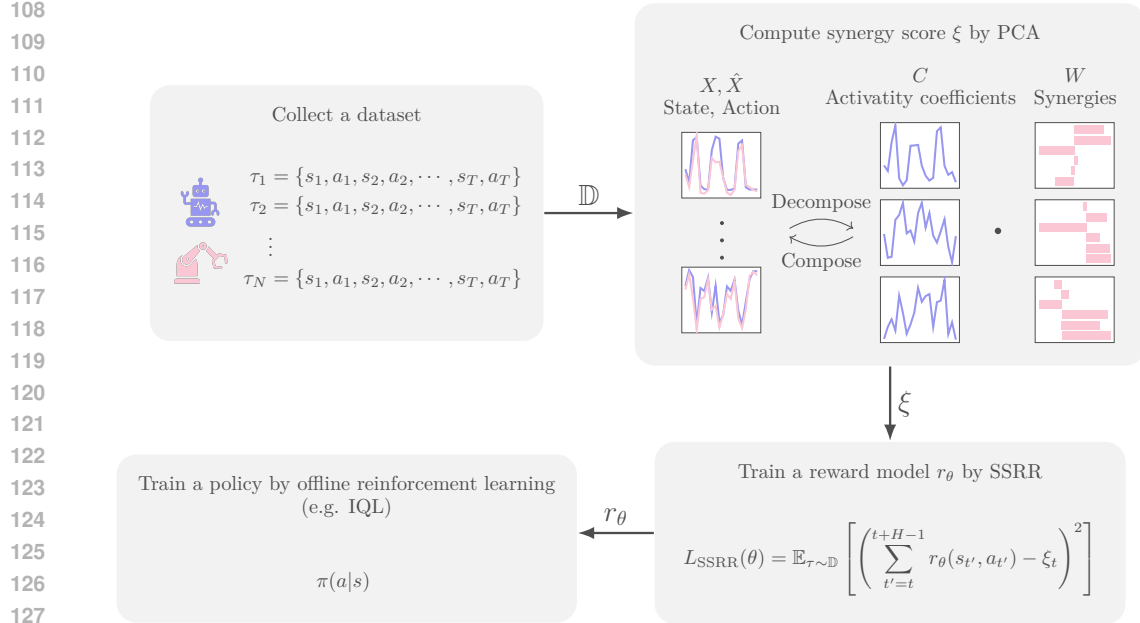


Figure 1: Overview of SynIL. Synergy scores ξ are calculated from trajectories of a dataset \mathbb{D} (top-right), and a reward model is trained based on the synergy scores (bottom-right). Offline reinforcement learning is finally performed with rewards computed by the trained reward model (bottom-left).

defined by the following equation:

$$\mathbf{X}_t = \mathbf{W}_t^n \mathbf{C}_t^n + \bar{\mathbf{X}}_t + \text{residuals}, \quad (1)$$

where

$$\mathbf{X}_t = \begin{bmatrix} \mathbf{s}_t & \mathbf{s}_{t+1} & \cdots & \mathbf{s}_{t+(H-1)} \\ \mathbf{a}_t & \mathbf{a}_{t+1} & \cdots & \mathbf{a}_{t+(H-1)} \end{bmatrix}. \quad (2)$$

Here, data \mathbf{X}_t is the H -step trajectory of the state-actions cut from a demonstration, n is the number of spatial synergies, \mathbf{W}_t^n corresponds to the spatial synergies, and \mathbf{C}_t^n represents the activation level of the spatial synergies at each time. Also, $\bar{\mathbf{X}}_t$ indicates the row-wise mean values of \mathbf{X}_t . Equation 1 is computed by the principle component analysis (PCA), which computes \mathbf{W}_t^n and \mathbf{C}_t^n to minimize residuals. It can be said that the smaller the residuals are, the more significant the spatial synergies are manifest.

Based on Equation 1, the degree of synergy manifestation is quantified by a metric denoted by $(R^2)_t^n$, as follows:

$$(R^2)_t^n = 1 - \frac{\|\mathbf{X}_t - \bar{\mathbf{X}}_t - \mathbf{W}_t^n \mathbf{C}_t^n\|_F^2}{\|\mathbf{X}_t - \bar{\mathbf{X}}_t\|_F^2}. \quad (3)$$

Here, $\|\bullet\|_F$ represents the Frobenius norm. $(R^2)_t^n$ ranges from 0 to 1, and the larger the value of $(R^2)_t^n$ is, the higher the degree of synergy manifestation is.

However, the above $(R^2)_t^n$ metric varies depending on the number of synergies, n . Instead, this study uses the synergy level (Chai & Hayashibe, 2020), ζ_t , as an indicator of synergy manifestation independent of n , defined as follows:

$$\zeta_t = \frac{1}{N} \sum_{n=1}^N (R^2)_t^n, \quad (4)$$

where N indicates the degree of freedom of demonstrations. Note that $(R^2)_t^0 \equiv 0$, and $(R^2)_t^n \equiv 1$ for $n \geq N$, so the synergy score sums up with a range of $n = 1, \dots, N$.

In addition to the synergy level, ζ , we also consider the standard deviation of the demonstrations. This is because only use of the synergy level can lead to the problem of high evaluation of demonstrations with overly simplistic states and actions, such as a demonstration with no movement. It is undesirable in many motor tasks to highly valuing the movement of doing nothing. To address this issue, we use the following synergy score, ξ_t , that considers the standard deviation of state-action pairs as well as the synergy level:

$$\xi_t = \log \zeta_t + \eta \log \bar{\sigma}_t. \quad (5)$$

Here, $\bar{\sigma}_t$ is the mean value of the standard deviation of columns of \mathbf{X}_t . η is a hyperparameter to weight the standard deviation of the demonstration action.

2.3 Reward Estimation Based on Synergy

As shown in Equations 2–5, the synergy score ξ_t are computed from an H -step state-action trajectory. However, to use offline reinforcement learning, we need rewards associated with each state-action pair. To bridge this gap, we consider that the synergy score ξ_t , computed from an H -step state-action trajectory, approximates the total reward of the H -step state-action trajectory, as expressed in the following equation:

$$\xi_t \approx \sum_{t'=t}^{t+H-1} r_{t'}, \quad (6)$$

where $r_{t'}$ indicates a reward at time t' .

To acquire such rewards, we construct a reward model that maps a state-action pair to a reward value. Let a reward model be $r_\theta(\mathbf{s}_t, \mathbf{a}_t)$, where \mathbf{s}_t and \mathbf{a}_t denote a state and action at time t , respectively, and θ indicates the parameters of the reward model. The reward model is trained to minimize the following loss function:

$$L_{SSRR}(\theta) = \mathbb{E}_{\tau \sim \mathcal{D}} \left[\left(\sum_{t'=t}^{t+H-1} r_\theta(s_{t'}, a_{t'}) - \xi_t \right)^2 \right]. \quad (7)$$

Minimizing this loss function gets the reward model to generate rewards which satisfy Equation 6. This method is similar to the self-supervised reward regression (SSRR) (Chen et al., 2021), with the noise level replaced with the synergy score.

3 Evaluation

3.1 Dataset

To evaluate the proposed method, we used D4RL(Fu et al., 2020) and Robomimic(Mandlekar et al., 2021) datasets as demonstration datasets. D4RL is a dataset for imitation learning and offline reinforcement learning, offering datasets for various environments and problem settings. Among various tasks, we specifically used the locomotion dataset as a typical motor task. The locomotion dataset includes three types of robots: Hopper, HalfCheetah, and Walker2d. D4RL contains several types of datasets with varying quality in terms of the performance and rewards of the expert policies. The “medium” dataset contains demonstrations with moderate quality, the “medium-replay” dataset is a mixed dataset of low to moderate quality demonstrations, and the “medium-expert” dataset is also a mix dataset of medium to high quality demonstrations. Each dataset was generated by policies at different training stages of the reinforcement learning algorithm, soft actor-critic(Haarhoja et al., 2018). These datasets are suitable for our problem setting, as they involve varying-quality demonstrations. It is known that typical imitation learning, such as behavior cloning (BC), decreases the performance with these types of datasets.

Robomimic, on the other hand, is an imitation-learning dataset with situations similar to real-world applications. The dataset is composed of demonstrations in various tasks with a seven degrees-of-freedom manipulator with a gripper. The demonstrations have been acquired by human operators with tele-operation in a simulator. Among datasets in

Robomimic, we used the multi-human (MH) datasets of Can and Square tasks. The MH dataset contains a total of 300 demonstrations, with 50 collected from each of six individuals with varying levels of expertise.

3.2 Problem Setting

In the evaluation, we assume a dataset can be expressed as follows:

$$\mathbb{D} = \{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) | t = 0, \dots, T - 1\}, \quad (8)$$

where \mathbf{s}_t and \mathbf{a}_t indicates a state and action at time t . We assume that the dataset does not have rewards or any sort of quality values for each tuple of $(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})$. Although the D4RL datasets contains rewards, we excluded them to emulate the situations where only demonstrations were provided. Thus, reinforcement learning methods cannot be applied unless we compute reward values.

We compared the proposed method with BC, a common imitation learning method. Additionally, for reference, we also compared them with the implicit Q-learning (IQL)(Kostrikov et al., 2021) with datasets containing the true rewards. This comparison was conducted to discuss how our method, which uses estimated rewards, is comparable with offline reinforcement learning with true rewards. The algorithm for SynIL is presented in Algorithm 1. Detail hyperparameters are provided in Appendix.

Algorithm 1 SynIL

Require: dataset \mathbb{D} of state-action trajectories, as defined in Equation 8

Ensure: a policy π and a reward model r_θ

 for each trajectory $\tau \sim \mathbb{D}$ do

 Compute ζ_t and $\bar{\sigma}_t$ from H -step trajectories cut from τ by Equation=4

 Compute ξ_t by Equation 5

 end for

 Initialize reward function parameters θ

 for each epoch do

 for each trajectory τ in \mathbb{D} do

 Compute loss $\mathcal{L}_{\text{SSRR}}(\theta)$ by Equation 7

 Update θ to minimize $\mathcal{L}_{\text{SSRR}}(\theta)$ through gradient descent

 end for

 end for

 Train a policy π by offline reinforcement learning with the trained reward model r_θ

3.3 Results

3.3.1 Correlation between Estimated Reward and Ground-Truth Reward

First, we investigated the relationship between the synergy-based rewards and ground-truth rewards in D4RL datasets. Figure 2 shows scatter plots of the rewards estimated by the SSRR (vertical axis, cumulative reward over 200 steps) and true rewards (horizontal axis, cumulative reward over 200 steps). Additionally, Table 1 presents the correlation coefficients between the rewards estimated by the SSRR and true rewards.

As a result, the correlation coefficients were more than 0.8 for all datasets except for the walker2d-medium-v2 and hopper-medium-replay-v2. Particularly, the correlation coefficients were more than 0.86 in all medium-expert datasets.

3.3.2 Performance Comparison

We evaluated the performance of imitation learning of SynIL (the proposed method), BC, and IQL in D4RL datasets. Table 2 shows the performance, which was computed as the average sum of rewards over 10 episodes. The values were normalized according to Fu et al. (2020).

270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323

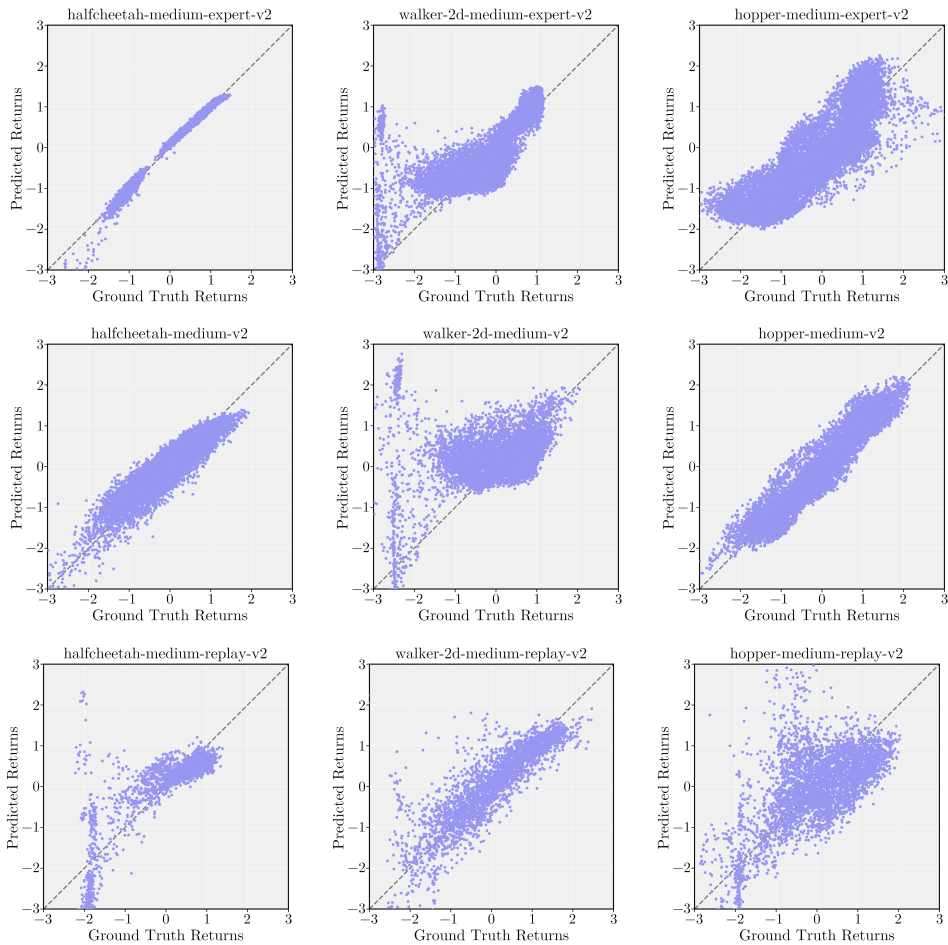


Figure 2: Comparison of synergy-based/true rewards for different datasets in D4RL. The vertical axes correspond to the cumulative rewards over 200 steps estimated by SSRR, and the horizontal axes correspond to the true cumulative rewards over 200 steps.

Table 1: Correlation coefficients between the synergy-based and true rewards in D4RL

Dataset	Correlation coefficient
halfcheetah-medium-expert-v2	0.9959
walker2d-medium-expert-v2	0.8655
hopper-medium-expert-v2	0.8718
halfcheetah-medium-v2	0.9482
walker2d-medium-v2	0.5809
hopper-medium-v2	0.9526
halfcheetah-medium-replay-v2	0.8039
walker2d-medium-replay-v2	0.8263
hopper-medium-replay-v2	0.5314

Table 2: Performance of each algorithm in D4RL (Average sum of rewards over 10 episodes)

Dataset	BC	SynIL (Ours)	IQL
halfcheetah-medium-v2	42.7 ± 0.3	45.4 ± 0.2	47.4 ± 0.2
hopper-medium-v2	54.4 ± 1.7	46.6 ± 3.8	54.4 ± 3.5
walker2d-medium-v2	71.7 ± 6.3	72.0 ± 5.4	80.0 ± 4.9
halfcheetah-medium-replay-v2	37.1 ± 2.4	40.6 ± 1.4	44.3 ± 0.5
hopper-medium-replay-v2	14.4 ± 2.4	99.9 ± 0.5	66.5 ± 8.1
walker2d-medium-replay-v2	15.5 ± 6.4	74.5 ± 3.0	69.5 ± 7.9
halfcheetah-medium-expert-v2	60.9 ± 5.0	90.9 ± 1.3	87.0 ± 3.1
hopper-medium-expert-v2	51.4 ± 1.8	80.9 ± 8.0	2.2 ± 0.0
walker2d-medium-expert-v2	80.0 ± 5.8	108.9 ± 0.3	110.7 ± 0.5

Table 3: Performance comparison in Robomimic (average success rate, 3 seeds × 50 rollouts)

Dataset	BC	SynIL (Ours)	IQL (sparse rewards)
Can (MH)	56.7 ± 10.0	94.0 ± 3.3	67.3 ± 5.0
Square (MH)	24.0 ± 4.3	49.3 ± 8.2	16.7 ± 5.2

As a result, according to Table 2, SynIL outperformed BC on most datasets and was comparable to IQL with true rewards. This is particularly evident in datasets of varying quality such as medium-replay and medium-expert. On the other hand, BC, which performs imitation by supervised learning, shows decreased performance on datasets with diverse quality like medium-expert.

We then evaluated the performance of imitation learning in Robomimic environments. Table 3 presents the performance, which was evaluated as the mean success rates for 50 trials across three different initial values (i.e. a total of 150 trials). Table 4 shows the execution time in seconds. Note that, since the Robomimic datasets did not contain reward data, we used sparse rewards, where a reward of 1 was applied at the final steps of successful episodes, and 0 was applied otherwise.

As a result, according to Table 3, SynIL significantly outperformed both BC and even IQL. Additionally, from Table 4, it is evident that SynIL had a shorter execution time than both BC and IQL.

Learning curves in D4RL and Robomimic are shown in Figures 3 and 4.

4 Discussion

4.1 Effectiveness of Synergy Rewards in Imitation Learning

The results suggest that the synergy-based rewards, calculated from the degree of synergy manifestation, can be used as surrogate rewards quantifying the quality of demonstrations. The synergy-based rewards have high correlation with the ground-truth rewards, as shown in Tables 1. Also, the use of the synergy-based rewards resulted in high performance in imitation learning, as shown in Tables 2, and 3. Those results demonstrate the effectiveness of the proposed method, that is, the use of synergy for evaluating the quality of demonstrations.

Table 4: Execution time in Robomimic (seconds)

Dataset	BC	SynIL (Ours)	IQL (sparse rewards)
Can (MH)	2.1 ± 0.2	1.3 ± 0.0	2.2 ± 0.2
Square (MH)	2.9 ± 0.0	2.0 ± 0.1	3.4 ± 0.2

378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431

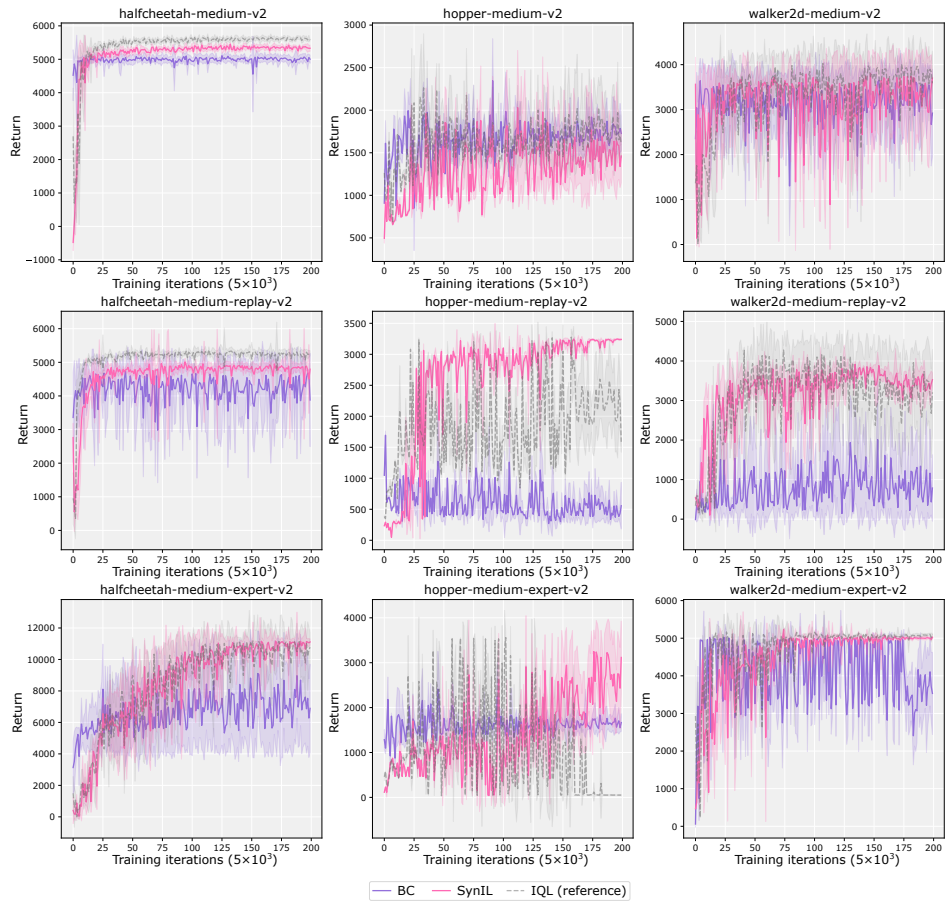


Figure 3: Learning curves of all environments.

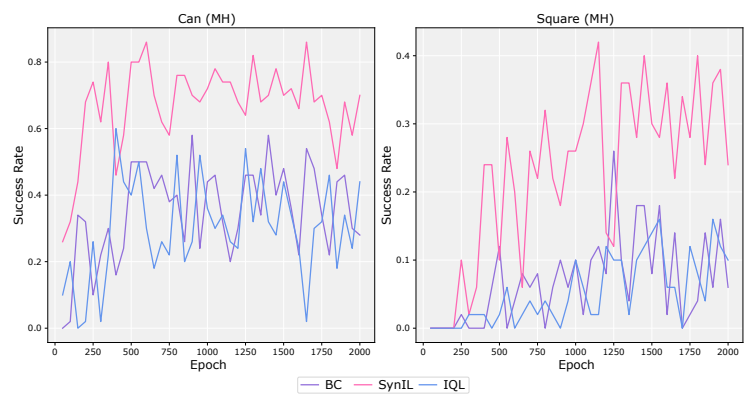


Figure 4: Learning curves of Robomimic.

432 The use of synergy-based rewards even outperformed IQL with ground-truth rewards in
433 some cases of D4RL and in Robomimic. This is because the ground-truth rewards we
434 used in Robomimic were sparse, which was hard to train a well-performing policy. The
435 synergy-based rewards, in contrast, could densely provide more information and function as
436 an important guide for successful completion in the task. This also highlights the advantage
437 of the synergy-based rewards that densely evaluate movement coordination.

439 4.2 Relationship Between Motor Learning and Synergy

440 The results shown in Tables 2 and 3 suggest that the synergy-based reward was effective
441 in improving the performance of imitation learning. Here, although it has not fully been
442 proven, we discuss why the synergy worked well in D4RL and Robomimic.

443 The success of the synergy-based method in D4RL would be caused by the nature of rein-
444 forcement learning. The demonstrations in D4RL datasets have been generated by neural-
445 network policies trained by deep reinforcement learning. In reinforcement learning, an agent
446 initially takes random state-action pairs, but as learning progresses, it takes more regular
447 states and actions. During this learning process, the dimensionality of the state-action dis-
448 tribution decreases. This links to the increase of synergy manifestation. Therefore, learning
449 state-action relationship where synergy is highly manifested will result in the selective learn-
450 ing of more proficient behaviors. This would be a reason why the proposed method resulted
451 in a high performance in the D4RL datasets. A similar phenomenon has been reported in
452 Chai & Hayashibe (2020); the degree of synergy manifestation increases as deep reinforce-
453 ment learning progresses. In future work, these findings may lead to a new reinforcement
454 learning framework that prefers to learn coordinated movements to accelerate motor learn-
455 ing. Moreover, previous works have reported that synergies obtained in some tasks can
456 generalize to new tasks (Al Borno et al., 2020; Kutsuzawa & Hayashibe, 2022). These find-
457 ings also suggest that the combination of synergy and reinforcement learning can improve
458 the ability of motor learning.

459 On the other hand, the reason of the success of the synergy-based method in Robomimic
460 would relate to the nature of the human motor learning. It is known that many human
461 movements become more regular, stable, and low-dimensional as learning progresses (Davids
462 et al., 2012). Even in unpredictable tasks, diverse actions are initially taken, but ultimately
463 they converge to typical patterns of adaptive responses due to the principle of optimality
464 (Braun et al., 2009). In such cases, it is considered that synergy in state actions also
465 manifests.

466 In addition to the above, the effectiveness of synergy can also be considered from the perspec-
467 tive of computational neuroscience. A recent theory in computational neuroscience argues
468 that human’s central nervous systems perform recognition and action generation so as to
469 minimize the variational free energy (Friston & Stephan, 2007). Minimizing the variational
470 free energy over time requires to decrease the amount of uncertainty, resulting in reduction
471 of the entropy of the sensory signals (Friston, 2012). Reduction in the entropy will bring
472 about more coordinated distributions of the sensory signals and actions, resulting in the
473 increase of the degree of synergy manifestation. Therefore, it is possible that as a human
474 becomes more proficient in a task, the variational free energy decreases more efficiently,
475 thereby the synergy is highly manifested.

476 From these perspectives and the results of this study, it is suggested that the synergy score
477 indicates the optimality of behaviors.

478 5 Conclusion

479 In this study, we aimed to address the issue in imitation learning that the performance
480 degrades due to varying qualities of demonstrations. We proposed a method of comput-
481 ing the quality of demonstrations based on synergy, enabling to use offline reinforcement
482 learning. We expected that synergy, a low-dimensional structure related to skill proficiency,
483 can be used to improve imitation learning performance. As a result, we reported that the
484 synergy-based imitation learning method outperformed the behavior cloning, a supervised
485

486 learning-based method. Although the relationship between synergy and demonstration qual-
487 ity has not yet been fully proven, we demonstrated the effectiveness of the use of synergy
488 across various datasets. We expect that our findings contribute to not only imitation learn-
489 ing but robotics and even computational neuroscience. Future work will focus on further
490 clarifying the scope of this method’s effectiveness and deepening the understanding of the
491 relationship between synergy and skill proficiency, potentially revealing the mechanisms of
492 skill acquisition in humans and other organisms.

493

494 References

495

496 Mazen Al Borno, Jennifer L. Hicks, and Scott L. Delp. The effects of motor modularity on
497 performance, learning and generalizability in upper-extremity reaching: a computational
498 analysis. *Journal of The Royal Society Interface*, 17(167):20200011, June 2020. doi:
499 10.1098/rsif.2020.0011.

500 Daniel A Braun, Ad Aertsen, Daniel M Wolpert, and Carsten Mehring. Learning optimal
501 adaptation strategies in unpredictable motor tasks. *Journal of Neuroscience*, 29(20):
502 6472–6478, 2009.

503 Jiazheng Chai and Mitsuhiro Hayashibe. Motor synergy development in high-performing
504 deep reinforcement learning algorithms. *IEEE Robotics and Automation Letters*, 5(2):
505 1271–1278, 2020.

506 Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstra-
507 tion via self-supervised reward regression. In *Conference on robot learning*, pp. 1262–1277.
508 PMLR, 2021.

509 Keith Davids, Duarte Araújo, Robert Hristovski, Pedro Passos, and Jia Yi Chow. Ecological
510 dynamics and motor learning design in sport. *Skill acquisition in sport: Research, theory
511 and practice*, 2:112–130, 2012.

512 Karl Friston. A free energy principle for biological systems. *Entropy*, 14(11):2100–2121,
513 2012.

514 Karl J Friston and Klaas E Stephan. Free-energy and the brain. *Synthese*, 159:417–458,
515 2007.

516 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets
517 for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.

518 Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile
519 manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*,
520 2024.

521 Reinhard Gentner, Susanne Gorges, David Weise, Kristin aufm Kampe, Mathias Buttmann,
522 and Joseph Classen. Encoding of motor skill in the corticomuscular system of musicians.
523 *Current Biology*, 20(20):1869–1874, 2010.

524 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-
525 policy maximum entropy deep reinforcement learning with a stochastic actor. In *Internat-
526 ional conference on machine learning*, pp. 1861–1870. PMLR, 2018.

527 Jihui Han, Jiazheng Chai, and Mitsuhiro Hayashibe. Synergy emergence in deep reinforce-
528 ment learning for full-dimensional arm manipulation. *IEEE Transactions on Medical
529 Robotics and Bionics*, 3(2):498–509, 2021.

530 Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. Vision-
531 language-action models for robotics: A review towards real-world applications. *IEEE
532 Access*, pp. 1–1, 2025. doi: 10.1109/ACCESS.2025.3609980.

533 Geon-Hyeong Kim, Seokin Seo, Jongmin Lee, Wonseok Jeon, HyeongJoo Hwang, Hongseok
534 Yang, and Kee-Eung Kim. Demodice: Offline imitation learning with supplementary
535 imperfect demonstrations. In *International Conference on Learning Representations*, 2022.
536
537
538
539

540 JW Kim, TZ Zhao, S Schmidgall, A Deguet, M Kobilarov, C Finn, and A Krieger. Surgical
541 robot transformer (srt): Imitation learning for surgical tasks. arxiv 2024. arXiv preprint
542 arXiv:2407.12998, 2024.

543

544 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
545 q-learning. arXiv preprint arXiv:2110.06169, 2021.

546 Kyo Kutsuzawa and Mitsuhiko Hayashibe. Motor synergy generalization framework for new
547 targets in multi-planar and multi-directional reaching task. Royal Society Open Science,
548 9(5):211721, 2022.

549

550 Mark L Latash, Nicholai A Bernstein, and Michael T Turvey. Dexterity and its development.
551 Psychology Press, 2014.

552 Ziniu Li, Tian Xu, Zeyu Qin, Yang Yu, and Zhi-Quan Luo. Imitation learning from im-
553 perfection: Theoretical justifications and algorithms. Advances in Neural Information
554 Processing Systems, 36, 2024.

555

556 Jinxin Liu, Lipeng Zu, Li He, and Donglin Wang. Clue: Calibrated latent guidance for
557 offline reinforcement learning. In Conference on Robot Learning, pp. 906–927. PMLR,
558 2023.

559 Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni,
560 Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in
561 learning from offline human demonstrations for robot manipulation. In arXiv preprint
562 arXiv:2108.03298, 2021.

563 Yunpeng Pan, Ching-An Cheng, Kamil Saigol, Keuntaek Lee, Xinyan Yan, Evangelos
564 Theodorou, and Byron Boots. Agile autonomous driving using end-to-end deep imitation
565 learning. arXiv preprint arXiv:1709.07174, 2017.

566

567 Michael T Turvey. Coordination. American psychologist, 45(8):938, 1990.

568 Haoran Xu, Xianyuan Zhan, Honglei Yin, and Huiling Qin. Discriminator-weighted of-
569 fline imitation learning from suboptimal demonstrations. In International Conference on
570 Machine Learning, pp. 24725–24742. PMLR, 2022.

571

572 Frank TJM Zaai, Kristin Daigle, Gerald L Gottlieb, and Esther Thelen. An unlearned
573 principle for controlling natural movements. Journal of Neurophysiology, 82(1):255–259,
574 1999.

575 Konrad Zolna, Alexander Novikov, Ksenia Konyushkova, Caglar Gulcehre, Ziyu Wang,
576 Yusuf Aytar, Misha Denil, Nando de Freitas, and Scott Reed. Offline learning from
577 demonstrations and unlabeled experience. arXiv preprint arXiv:2011.13885, 2020.

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Table 5: Hyperparameters for SSRR

	Hyperparameter	Value in D4RL	Value in Robomimic
RewardNet	Hidden layers	2	2
	Hidden Size	512	1024
	Activation	ReLU	ReLU
	Last Activation	Linear	Linear
	Epochs	500	10000
	Batch Size	64	64
	Learning Rate	1×10^{-4}	1×10^{-3}
	Weight Decay	1×10^{-2}	0
	Criterion	MSELoss	MSELoss
	Optimizer	Adam	Adam
Synergy Score	Horizon H	200	20
	Stride S	100	20
	η	0.05	0.05

Table 6: Hyperparameters for imitation learning

Hyperparameter	SynIL and IQL	BC
Hidden size	256	256
Hidden layers	2	2
Activation	ReLU	ReLU
Optimizer	Adam	Adam
Epochs	10^6	10^6
Evaluate period	5000	5000
Batch size	256	256
Learning rate	3×10^{-4}	3×10^{-4}
Discount	0.99	-
α	0.005	-
τ	0.7	-
β	3.0	-

A Appendix

A.1 Hyperparameters

The neural network and Synergy Score hyperparameters used in the SSRR are listed in Table 5. The input dimension of the reward estimation model (RewardNet) was equal to the sum of the state and action dimensions, and the output dimension was one-dimensional since it represents the reward. H and S indicate the horizon and stride; H -step trajectories are cut from a demonstration while shifting the start position by S -steps. Additionally, η was set through hand-tuning to achieve the highest correlation with the true reward.

Hyperparameters for SynIL, IQL, and BC are shown in Table 6. IQL becomes equivalent to the BC algorithm by setting the hyperparameter β to 0. Therefore, in D4RL, BC experiments were conducted by setting IQL’s β to 0. The hyperparameters were hand-tuned to achieve the highest performance for each algorithm.