

---

# COUNTERFACTUAL TRAINING: TEACHING MODELS PLAUSIBLE AND ACTIONABLE EXPLANATIONS

---

A PREPRINT

**Patrick Altmeyer** 

Faculty of Electrical Engineering, Mathematics and Computer Science  
Delft University of Technology

[p.altmeyer@tudelft.nl](mailto:p.altmeyer@tudelft.nl)

**Arie van Deursen**

Faculty of Electrical Engineering, Mathematics and Computer Science  
Delft University of Technology

**Cynthia C. S. Liem**

Faculty of Electrical Engineering, Mathematics and Computer Science  
Delft University of Technology

January 31, 2025

## ABSTRACT

Counterfactual Explanations (CE) have emerged as a popular tool to explain predictions made by opaque machine learning models: they explain how factual inputs need to change in order for some fitted model to produce some desired output. Much existing research has focused on identifying explanations that are not only valid but also deemed desirable with respect to the underlying data and stakeholder requirements. Recent work has shown that under this premise, the task of learning desirable explanations is effectively reassigned from the model itself to the (post-hoc) counterfactual explainer. Building on that work, we propose a novel model objective that leverages counterfactuals during the training phase (ad-hoc) in order to minimize the divergence between learned representations and desirable explanations. Through extensive experiments, we demonstrate that our proposed methodology facilitates training models that inherently deliver desirable explanations while maintaining high predictive performance.

**Keywords** Counterfactual Explanations • Explainable AI • Representation Learning

## 1 Introduction

Today's prominence of artificial intelligence (AI) has largely been driven by advances in **representation learning**: instead of relying on features and rules that are carefully hand-crafted by humans, modern AIs are tasked with learning these representations from scratch, guided by narrow objectives such as predictive accuracy (I. Goodfellow, Bengio, and Courville 2016). Modern advances in computing have made it possible to provide such AIs with ever greater degrees of freedom to achieve that task, which has often led them to outperform traditionally more parsimonious models. Unfortunately, in doing so they have also learned increasingly complex representations that we can no longer easily interpret.

This trend towards complexity for the sake of performance has come under serious scrutiny in recent years. At the very cusp of the deep learning revolution, I. J. Goodfellow, Shlens, and Szegedy (2014) showed that artificial neural

23 networks (ANN) are highly sensitive to adversarial examples (AE): counterfactuals of model inputs that yield vastly  
 24 different model predictions despite being semantically indifferent from their factual counterparts. Despite continued  
 25 research efforts towards making models more robust to such counterfactuals, the issue remains unresolved even for  
 26 models that are considered to small by today's standards (Kolter 2023). It seems that robust representation learning is  
 27 fundamentally at odds with large degrees of freedom.

28 Wilson (2020)

29 Rudin (2019)

## 30 2 Related Literature

### 31 2.1 Background on Counterfactual Explanations

32 (Wachter, Mittelstadt, and Russell 2017; Joshi et al. 2019; Altmeyer et al. 2024)

### 33 2.2 Learning Representations

34 For example, joint-energy models

### 35 2.3 Generalization and Robustness

36 Sauer and Geiger (2021) generate counterfactual images for MNIST and ImageNet through independent mechanisms  
 37 (IM): each IM learns class-conditional input distributions over a specific lower-dimensional, semantically meaningful  
 38 factor, such as *texture*, *shape* and *background*. They demonstrate that using these generated counterfactuals during  
 39 classifier training improves model robustness. Similarly, Abbasnejad et al. (2020) argue that counterfactuals represent  
 40 potentially useful training data in machine learning, especially in supervised settings where inputs may be reasonably  
 41 mapped to multiple outputs. They, too, demonstrate the augmenting the training data of image classifiers can improve  
 42 generalization.

43 Teney, Abbasnejad, and Hengel (2020) propose an approach using counterfactuals in training that does not rely on  
 44 data augmentation: they argue that counterfactual pairs typically already exist in training datasets. Specifically, their  
 45 approach relies on, firstly, identifying similar input samples with different annotations and, secondly, ensuring that the  
 46 gradient of the classifier aligns with the vector between pairs of counterfactual inputs using the cosine distance as a loss  
 47 function (referred to as *gradient supervision*) (*this might be useful for our task as well*). In the natural language pro-  
 48 cessing (NLP) domain, counterfactuals have similarly been used to improve models through data augmentation: Wu et  
 49 al. (2021), propose POLYJUICE, a general-purpose counterfactual generator for language models. They demonstrate  
 50 empirically that augmenting training data through POLYJUICE counterfactuals improves robustness in a number of  
 51 NLP tasks.

### 52 2.4 Link to Adversarial Training

53 Freiesleben (2022) propose two definitional differences between Adversarial Examples (AE) and Counterfactual Ex-  
 54 planations (CE): firstly, and more importantly according to the authors, the term AE implies missclassification, which  
 55 is not the case for CE (*this might be a useful notion for use to distinguish between adversarials and explanations*  
 56 *during training*); secondly, they argue that closeness plays a more critical role in the context of CE but confess that  
 57 even counterfactuals that are not close might be relevant explanations. Pawelczyk et al. (2022) show that CE and AE  
 58 are equivalent under certain conditions and derive upper bounds on the distances between them.

### 59 2.5 Closely Related

60 Guo, Nguyen, and Yadav (2023) are the first to propose end-to-end training pipeline that includes counterfactual ex-  
 61 planations as part of the training procedure. In particular, they propose a specific network architecture that includes  
 62 a predictor and CE generator network (*akin a GAN?*), where the parameters of the CE generator network are learn-  
 63 able. Counterfactuals are generated during each training iteration and fed back to the predictor network (*here we are*  
 64 *aligned*). In contrast, we impose no restrictions on the neural network architecture at all. (*to ensure the one-hot en-*  
 65 *coding of categorical features is maintained, they simple use softmax (might be interesting for CE.jl)*) Interestingly,  
 66 the authors find that their approach is sensitive to the choice of the loss function: only MSE seems to lead to good  
 67 performance. They also demonstrate theoretically, that the objective function is difficult to optimize due to divergent  
 68 gradients and suffers from poor adversarial robustness. (*because partial gradients with respect to the classification*  
 69 *loss component and the counterfactual validity component point in opposite directions*). To mitigate these issues,  
 70 the authors use block-wise gradient descent: they first update with respect to classification loss and then use a second  
 71 update with respect to the other loss components (*this might be useful for our task as well*). Ross, Lakkaraju, and  
 72 Bastani (2024) propose a way to train models that are guaranteed to provide recourse for individuals with high proba-  
 73 bility. The approach builds on adversarial training (*here we are aligned*), where in this context adversarial examples

- 74 are actively encouraged to exist, but only target attacks with respect to the positive class. The proposed method allows  
 75 for imposing a set of actionable recourse ex-ante: for example, users can impose mutability constraints for features  
 76 (*here we are aligned*). (*To solve their objective function more efficiently, they use a first-order Taylor approximation*  
 77 *to approximate the recourse loss component (might be applicable in our case)*)
- 78 Luu and Inoue (2023) introduce Counterfactual Adversarial Training (CAT) with intention of improving generalization  
 79 and robustness of language models. Specifically, they propose to proceed as follows: firstly, identify training samples  
 80 that are subject to high predictive uncertainty (entropy); secondly, generate counterfactual explanations for those  
 81 samples; and, finally, finetune the model on the augmented dataset that includes the generated counterfactuals.

82 **3 Counterfactual Training**

83 **4 Experiments**

84 **4.1 Experimental Setup**

85 **4.2 Experimental Results**

86 **5 Discussion**

87 **6 Conclusion**

88 **References**

- 89 Abbasnejad, Ehsan, Damien Teney, Amin Parvaneh, Javen Shi, and Anton van den Hengel. 2020. “Counterfactual  
 90 Vision and Language Learning.” In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
 91 (*CVPR*), 10041–51. <https://doi.org/10.1109/CVPR42600.2020.01006>.
- 92 Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia CS Liem. 2024. “Faithful Model Explanations  
 93 Through Energy-Constrained Conformal Counterfactuals.” In *Proceedings of the AAAI Conference on Artificial*  
 94 *Intelligence*, 38:10829–37. 10.
- 95 Freiesleben, Timo. 2022. “The Intriguing Relation Between Counterfactual Explanations and Adversarial Examples.”  
*Minds and Machines* 32 (1): 77–109.
- 96 Goodfellow, Ian J, Jonathon Shlens, and Christian Szegedy. 2014. “Explaining and Harnessing Adversarial Examples.”  
<https://arxiv.org/abs/1412.6572>.
- 97 Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- 98 Guo, Hangzhi, Thanh H. Nguyen, and Amulya Yadav. 2023. “CounterNet: End-to-End Training of Prediction Aware  
 99 Counterfactual Explanations.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery*  
 100 *and Data Mining*, 577–89. KDD ’23. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3580305.3599290>.
- 101 Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vigitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. “Towards Realistic  
 102 Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems.” <https://arxiv.org/abs/1907.09615>.
- 103 Kolter, Zico. 2023. “Keynote Addresses: SaTML 2023 .” In *2023 IEEE Conference on Secure and Trustworthy*  
*Machine Learning (SaTML)*, xvi–. Los Alamitos, CA, USA: IEEE Computer Society. <https://doi.org/10.1109/SaTML54575.2023.00009>.
- 104 Luu, Hoai Linh, and Naoya Inoue. 2023. “Counterfactual Adversarial Training for Improving Robustness of Pre-  
 105 Trained Language Models.” In *Proceedings of the 37th Pacific Asia Conference on Language, Information and*  
 106 *Computation*, 881–88.
- 107 Pawelczyk, Martin, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, and Himabindu Lakkaraju. 2022. “Exploring  
 108 Counterfactual Explanations Through the Lens of Adversarial Examples: A Theoretical and Empirical Analysis.”  
 109 In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, edited by Gustau  
 110 Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, 151:4574–94. Proceedings of Machine Learning Research.  
 111 PMLR. <https://proceedings.mlr.press/v151/pawelczyk22a.html>.
- 112 Ross, Alexis, Himabindu Lakkaraju, and Osbert Bastani. 2024. “Learning Models for Actionable Recourse.” In  
 113 *Proceedings of the 35th International Conference on Neural Information Processing Systems*. NIPS ’21. Red  
 114 Hook, NY, USA: Curran Associates Inc.
- 115 Rudin, Cynthia. 2019. “Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use  
 116 Interpretable Models Instead.” *Nature Machine Intelligence* 1 (5): 206–15. <https://doi.org/10.1038/s42256-019-0048-x>.
- 117 Sauer, Axel, and Andreas Geiger. 2021. “Counterfactual Generative Networks.” <https://arxiv.org/abs/2101.06046>.

- 125 Tenev, Damien, Ehsan Abbasnedjad, and Anton van den Hengel. 2020. “Learning What Makes a Difference from  
126 Counterfactual Examples and Gradient Supervision.” In *Computer Vision–ECCV 2020: 16th European Confer-  
127 ence, Glasgow, UK, August 23–28, 2020, Proceedings, Part x 16*, 580–99. Springer.  
128 Wachter, Sandra, Brent Mittelstadt, and Chris Russell. 2017. “Counterfactual Explanations Without Opening the Black  
129 Box: Automated Decisions and the GDPR.” *Harv. JL & Tech.* 31: 841. <https://doi.org/10.2139/ssrn.3063289>.  
130 Wilson, Andrew Gordon. 2020. “The Case for Bayesian Deep Learning.” <https://arxiv.org/abs/2001.10995>.  
131 Wu, Tongshuang, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. “Polyjuice: Generating Counterfactuals  
132 for Explaining, Evaluating, and Improving Models.” In *Proceedings of the 59th Annual Meeting of the Associa-  
133 tion for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing  
134 (Volume 1: Long Papers)*, edited by Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, 6707–23. Online:  
135 Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-long.523>.

136 **A Training Details**

137 **A.1 Initial Grid Search**

138 For the initial round of experiments we

139 **A.1.1 Generator Parameters**

140 The hyperparameter grids for the first investigation of the effect of generator parameters are shown in Parameters [A.1](#)  
141 and Parameters [A.2](#).

142 **Parameters A.1 (Training Phase).**

- 143 • Generator Parameters:
  - 144 –  $\lambda_{\text{cost}}$ : 0.0, 0.001, 0.1
  - 145 –  $\lambda_{\text{div}}$ : 0.01, 0.05, 0.1, 0.5, 1.0, 5.0, 10.0, 15.0
  - 146 – Learning Rate: 1.0
  - 147 – Maximum Iterations: 20, 50, 100
  - 148 – Optimizerimizer: sgd
- 149 • Generator: `ecco`, `generic`, `omni`, `revise`
- 150 • Training Parameters:
  - 151 – Objective: `full`, `vanilla`

152 **Parameters A.2 (Evaluation Phase).**

- 153 • Counterfactual Parameters:
  - 154 – Convergence: `max_iter`
  - 155 – Maximum Iterations: 100
  - 156 – No. Individuals: 100
  - 157 – No. Runs: 5
- 158 • Generator Parameters:
  - 159 –  $\lambda_{\text{cost}}$ : 0.0
  - 160 –  $\lambda_{\text{div}}$ : 0.1, 0.5, 1.0, 5.0, 10.0, 20.0
  - 161 – Learning Rate: 1.0
  - 162 – Maximum Iterations: 50
  - 163 – Optimizerimizer: sgd

164 **A.1.1.1 Linearly Separable**

- 165 • **Energy Penalty** (Table [A1](#)): *ECCo* generally does yield better results than *Vanilla* for higher choices of the  
166 energy penalty (10,15) during training. *Generic* performs poorly accross the board. *Omni* seems to have an  
167 anchoring effect, in that it never performs terribly but also never as good as the best *ECCo* results. *REVISE*  
168 performs poorly across the board.
- 169 • **Cost** (Table [A2](#)): Results for all generators (except *Omni*) are quite bad, which can likely be attributed to  
170 extremely bad results for some choices of the **Energy Penalty** (results here are averaged). For *ECCo* and  
171 *Generic*, higher cost values generally lead to worse results.
- 172 • **Maximum Iterations**: No clear patterns recognizable, so it seems that smaller choices are ok.
- 173 • **Validity**: *ECCo* almost always valid except for very low values during training and high values at evaluation  
174 time. *Generic* often has poor validity.
- 175 • **Accuracy**: Seems largely unaffected.

Table A1: Results for Linearly Separable data by energy penalty.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
full	0.01	<i>ECCo</i>	$-9.91 \cdot 10^{11}$	$2.25 \cdot 10^{12}$
full	0.01	<i>Generic</i>	$-5.71 \cdot 10^{17}$	$1.3 \cdot 10^{18}$
<b>full</b>	<b>0.01</b>	<b>Omniscient</b>	<b>-2.54</b>	<b>0.116</b>
full	0.01	<i>REVISE</i>	-15.6	13.2

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
vanilla	0.01	<i>ECCo</i>	-4.28	3.52
vanilla	0.01	<i>Generic</i>	-4.45	3.47
vanilla	0.01	<i>Omniscient</i>	-5.12	4.46
vanilla	0.01	<i>REVISE</i>	-4.91	4.24
full	0.05	<i>ECCo</i>	$-5.63 \cdot 10^5$	$1.28 \cdot 10^6$
full	0.05	<i>Generic</i>	$-8.35 \cdot 10^{17}$	$1.9 \cdot 10^{18}$
<b>full</b>	<b>0.05</b>	<b>Omniscient</b>	<b>-2.53</b>	<b>0.114</b>
full	0.05	<i>REVISE</i>	-15	12.6
vanilla	0.05	<i>ECCo</i>	-4.4	3.66
vanilla	0.05	<i>Generic</i>	-4.38	3.48
vanilla	0.05	<i>Omniscient</i>	-5.25	4.62
vanilla	0.05	<i>REVISE</i>	-4.94	4.22
full	0.1	<i>ECCo</i>	$-6.74 \cdot 10^5$	$1.53 \cdot 10^6$
full	0.1	<i>Generic</i>	$-1.72 \cdot 10^{11}$	$3.9 \cdot 10^{11}$
<b>full</b>	<b>0.1</b>	<b>Omniscient</b>	<b>-2.56</b>	<b>0.124</b>
full	0.1	<i>REVISE</i>	-15.6	13.2
vanilla	0.1	<i>ECCo</i>	-4.28	3.52
vanilla	0.1	<i>Generic</i>	-4.45	3.48
vanilla	0.1	<i>Omniscient</i>	-5.12	4.46
vanilla	0.1	<i>REVISE</i>	-4.91	4.25
full	0.5	<i>ECCo</i>	-11.8	9.83
full	0.5	<i>Generic</i>	$-1.06 \cdot 10^{18}$	$2.42 \cdot 10^{18}$
<b>full</b>	<b>0.5</b>	<b>Omniscient</b>	<b>-2.54</b>	<b>0.123</b>
full	0.5	<i>REVISE</i>	-15	12.6
vanilla	0.5	<i>ECCo</i>	-4.4	3.65
vanilla	0.5	<i>Generic</i>	-4.38	3.48
vanilla	0.5	<i>Omniscient</i>	-5.25	4.61
vanilla	0.5	<i>REVISE</i>	-4.95	4.22
full	1	<i>ECCo</i>	-11.5	11.1
full	1	<i>Generic</i>	$-1.71 \cdot 10^{11}$	$3.88 \cdot 10^{11}$
<b>full</b>	<b>1</b>	<b>Omniscient</b>	<b>-2.59</b>	<b>0.117</b>
full	1	<i>REVISE</i>	-15.7	13.3
vanilla	1	<i>ECCo</i>	-4.28	3.51
vanilla	1	<i>Generic</i>	-4.44	3.47
vanilla	1	<i>Omniscient</i>	-5.11	4.46
vanilla	1	<i>REVISE</i>	-4.91	4.25
full	5	<i>ECCo</i>	-3.99	3.12
full	5	<i>Generic</i>	$-4.88 \cdot 10^{17}$	$1.11 \cdot 10^{18}$
<b>full</b>	<b>5</b>	<b>Omniscient</b>	<b>-2.53</b>	<b>0.117</b>
full	5	<i>REVISE</i>	-14.6	12.1
vanilla	5	<i>ECCo</i>	-4.4	3.65
vanilla	5	<i>Generic</i>	-4.38	3.48
vanilla	5	<i>Omniscient</i>	-5.25	4.61
vanilla	5	<i>REVISE</i>	-4.95	4.22
<b>full</b>	<b>10</b>	<b>ECCo</b>	<b>-2.31</b>	<b>0.735</b>
full	10	<i>Generic</i>	$-1.7 \cdot 10^{11}$	$3.86 \cdot 10^{11}$
full	10	<i>Omniscient</i>	-2.53	0.117
full	10	<i>REVISE</i>	-15.5	13
vanilla	10	<i>ECCo</i>	-4.28	3.51
vanilla	10	<i>Generic</i>	-4.44	3.47
vanilla	10	<i>Omniscient</i>	-5.12	4.46
vanilla	10	<i>REVISE</i>	-4.91	4.24
<b>full</b>	<b>15</b>	<b>ECCo</b>	<b>-2.01</b>	<b>0.488</b>
full	15	<i>Generic</i>	$-4.91 \cdot 10^{17}$	$1.12 \cdot 10^{18}$
full	15	<i>Omniscient</i>	-2.53	0.116

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
full	15	<i>REVISE</i>	-14.4	11.7
vanilla	15	<i>ECCo</i>	-4.4	3.65
vanilla	15	<i>Generic</i>	-4.38	3.48
vanilla	15	<i>Omniscient</i>	-5.25	4.6
vanilla	15	<i>REVISE</i>	-4.95	4.23

Table A2: Results for Linearly Separable data by cost penalty.

Objective	$\lambda_{\text{cost}}(\text{train})$	Generator	Value	Std
full	0	<i>ECCo</i>	$-5.32 \cdot 10^3$	$1.21 \cdot 10^4$
full	0	<i>Generic</i>	$-1.03 \cdot 10^{18}$	$2.34 \cdot 10^{18}$
<b>full</b>	<b>0</b>	<b>Omniscient</b>	<b>-2.64</b>	<b>0.125</b>
full	0	<i>REVISE</i>	-15.4	12.9
vanilla	0	<i>ECCo</i>	-4.34	3.58
vanilla	0	<i>Generic</i>	-4.41	3.48
vanilla	0	<i>Omniscient</i>	-5.18	4.54
vanilla	0	<i>REVISE</i>	-4.93	4.23
full	0.001	<i>ECCo</i>	-362	811
full	0.001	<i>Generic</i>	$-2.65 \cdot 10^{17}$	$6.03 \cdot 10^{17}$
<b>full</b>	<b>0.001</b>	<b>Omniscient</b>	<b>-2.49</b>	<b>0.115</b>
full	0.001	<i>REVISE</i>	-15.5	13
vanilla	0.001	<i>ECCo</i>	-4.34	3.58
vanilla	0.001	<i>Generic</i>	-4.41	3.48
vanilla	0.001	<i>Omniscient</i>	-5.18	4.53
vanilla	0.001	<i>REVISE</i>	-4.93	4.23
full	0.1	<i>ECCo</i>	$-3.72 \cdot 10^{11}$	$8.46 \cdot 10^{11}$
full	0.1	<i>Generic</i>	$-4.49 \cdot 10^{14}$	$1.02 \cdot 10^{15}$
<b>full</b>	<b>0.1</b>	<b>Omniscient</b>	<b>-2.5</b>	<b>0.112</b>
full	0.1	<i>REVISE</i>	-14.6	12.2
vanilla	0.1	<i>ECCo</i>	-4.34	3.58
vanilla	0.1	<i>Generic</i>	-4.41	3.48
vanilla	0.1	<i>Omniscient</i>	-5.18	4.54
vanilla	0.1	<i>REVISE</i>	-4.93	4.24

## 176 A.1.1.2 Moons

- 177 • **Energy Penalty** (Table A3): *ECCo* consistently yields better results than *Vanilla*, except for very low choices  
 178 of the energy penalty during training for which it performs abysmal. *Generic* performs quite badly across  
 179 the board for high enough choices of the energy penalty at evaluation time. *Omni* has small positive effect.  
 180 *REVISE* performs poorly across the board.
- 181 • **Cost (distance penalty)**: *Generic* generally does better for higher values, while *ECCo* does better for lower  
 182 values.
- 183 • **Maximum Iterations**: No clear patterns recognizable, so it seems that smaller choices are ok.
- 184 • **Validity**: *ECCo* generally achieves full validity except for very low choices the energy penalty during training  
 185 and high choices at evaluation time. *Generic* performs poorly for high choices of the energy penalty during  
 186 evaluation.
- 187 • **Accuracy**: Largely unaffected although *ECCo* suffers a bit for very low choices the energy penalty during  
 188 training. *REVISE* suffers a lot in general (around 10 percentage points).

Table A3: Results for Moons data by energy penalty.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
full	0.01	<i>ECCo</i>	$-2.8 \cdot 10^{22}$	$6.39 \cdot 10^{22}$
full	0.01	<i>Generic</i>	$-4.89 \cdot 10^{30}$	$1.11 \cdot 10^{31}$
<b>full</b>	<b>0.01</b>	<b>Omniscient</b>	<b>-4.74</b>	<b>5.08</b>
full	0.01	<i>REVISE</i>	-572	$1.25 \cdot 10^3$
vanilla	0.01	<i>ECCo</i>	-15.5	17.3
vanilla	0.01	<i>Generic</i>	-10.9	11.9
vanilla	0.01	<i>Omniscient</i>	-12.7	14.4
vanilla	0.01	<i>REVISE</i>	-11.2	13
full	0.05	<i>ECCo</i>	$-1.55 \cdot 10^{16}$	$3.52 \cdot 10^{16}$
full	0.05	<i>Generic</i>	$-2.22 \cdot 10^{20}$	$5 \cdot 10^{20}$
<b>full</b>	<b>0.05</b>	<b>Omniscient</b>	<b>-4.41</b>	<b>4.48</b>
full	0.05	<i>REVISE</i>	$-1.04 \cdot 10^3$	$2.3 \cdot 10^3$
vanilla	0.05	<i>ECCo</i>	-15.5	17.2
vanilla	0.05	<i>Generic</i>	-11.7	12.8
vanilla	0.05	<i>Omniscient</i>	-12.4	14.1
vanilla	0.05	<i>REVISE</i>	-11.3	13.1
full	0.1	<i>ECCo</i>	$-3.41 \cdot 10^3$	$7.73 \cdot 10^3$
full	0.1	<i>Generic</i>	$-5.22 \cdot 10^{30}$	$1.19 \cdot 10^{31}$
<b>full</b>	<b>0.1</b>	<b>Omniscient</b>	<b>-4.78</b>	<b>5.12</b>
full	0.1	<i>REVISE</i>	-288	594
vanilla	0.1	<i>ECCo</i>	-15.5	17.2
vanilla	0.1	<i>Generic</i>	-10.9	11.9
vanilla	0.1	<i>Omniscient</i>	-12.7	14.4
vanilla	0.1	<i>REVISE</i>	-11.3	13.1
full	0.5	<i>ECCo</i>	-7.09	7.51
full	0.5	<i>Generic</i>	$-1.11 \cdot 10^{31}$	$2.53 \cdot 10^{31}$
<b>full</b>	<b>0.5</b>	<b>Omniscient</b>	<b>-4.58</b>	<b>4.83</b>
full	0.5	<i>REVISE</i>	$-1.19 \cdot 10^3$	$2.64 \cdot 10^3$
vanilla	0.5	<i>ECCo</i>	-15.5	17.2
vanilla	0.5	<i>Generic</i>	-11.7	12.8
vanilla	0.5	<i>Omniscient</i>	-12.4	14.1
vanilla	0.5	<i>REVISE</i>	-11.3	13.1
full	1	<i>ECCo</i>	-6.06	6.33
full	1	<i>Generic</i>	$-1.58 \cdot 10^{33}$	$3.59 \cdot 10^{33}$
<b>full</b>	<b>1</b>	<b>Omniscient</b>	<b>-4.66</b>	<b>4.89</b>
full	1	<i>REVISE</i>	$-1.16 \cdot 10^3$	$2.59 \cdot 10^3$
vanilla	1	<i>ECCo</i>	-15.5	17.3
vanilla	1	<i>Generic</i>	-10.9	11.9
vanilla	1	<i>Omniscient</i>	-12.7	14.4
vanilla	1	<i>REVISE</i>	-11.3	13.1
<b>full</b>	<b>5</b>	<b>ECCo</b>	<b>-2.57</b>	<b>2.07</b>
full	5	<i>Generic</i>	$-1.17 \cdot 10^{28}$	$2.66 \cdot 10^{28}$
full	5	<i>Omniscient</i>	-4.29	4.31
full	5	<i>REVISE</i>	-530	$1.16 \cdot 10^3$
vanilla	5	<i>ECCo</i>	-15.5	17.2
vanilla	5	<i>Generic</i>	-11.7	12.7
vanilla	5	<i>Omniscient</i>	-12.4	14.1
vanilla	5	<i>REVISE</i>	-11.3	13.1
<b>full</b>	<b>10</b>	<b>ECCo</b>	<b>-1.76</b>	<b>0.974</b>
full	10	<i>Generic</i>	$-1.54 \cdot 10^{33}$	$3.51 \cdot 10^{33}$
full	10	<i>Omniscient</i>	-4.44	4.56
full	10	<i>REVISE</i>	$-1.52 \cdot 10^3$	$3.4 \cdot 10^3$
vanilla	10	<i>ECCo</i>	-15.5	17.3

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
vanilla	10	<i>Generic</i>	-10.9	11.9
vanilla	10	<i>Omniscient</i>	-12.7	14.4
vanilla	10	<i>REVISE</i>	-11.3	13.1
<b>full</b>	<b>15</b>	<b>ECCo</b>	<b>-1.37</b>	<b>0.365</b>
full	15	<i>Generic</i>	$-5.32 \cdot 10^{28}$	$1.21 \cdot 10^{29}$
full	15	<i>Omniscient</i>	-4.34	4.38
full	15	<i>REVISE</i>	-473	$1.03 \cdot 10^3$
vanilla	15	<i>ECCo</i>	-15.5	17.2
vanilla	15	<i>Generic</i>	-11.7	12.8
vanilla	15	<i>Omniscient</i>	-12.4	14.1
vanilla	15	<i>REVISE</i>	-11.3	13.1

## 189 A.1.1.3 Circles

- 190 • **Energy Penalty** (Table A4): *ECCo* consistently yields better results than *Vanilla*, though primarily for low to  
 191 medium choices of the energy penalty ( $<=5$ ) during training. The same goes for *Generic*, which sometimes  
 192 outperforms *ECCo* (for small energy penalty at evaluation time). *Omni* does alright for lower energy penalty  
 193 at evaluation time, but loses out for higher choices. *REVISE* performs poorly across the board (except very  
 194 low choices at evaluation time).
- 195 • **Cost (distance penalty)**: *ECCo* and *Generic* generally achieve the best results when no cost penalty is used  
 196 during training. Both *Omni* and *REVISE* are largely unaffected.
- 197 • **Maximum Iterations**: *ECCo* consistently yields better results for higher numbers of iterations. *Generic*  
 198 generally does best for a medium number (50). *Omni* is sometimes invalid (???).
- 199 • **Validity**: *ECCo* tends to outperform its *Vanilla* counterpart, though primarily for low to medium choices of  
 200 the energy penalty ( $<=5$ ) during training and evaluation. *Vanilla* typically worse across the board.
- 201 • **Accuracy**: Mostly unaffected, but *REVISE* again consistently some deterioration and *ECCo* deteriorates for  
 202 high choices of energy penalty during training, reflecting other outcomes above.

Table A4: Results for Circles data by energy penalty.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
<b>full</b>	<b>0.01</b>	<b>ECCo</b>	<b>-1.26</b>	<b>0.423</b>
full	0.01	<i>Generic</i>	-1.49	0.71
full	0.01	<i>Omniscient</i>	-5.21	5.25
full	0.01	<i>REVISE</i>	$-2.71 \cdot 10^{26}$	$6.37 \cdot 10^{26}$
vanilla	0.01	<i>ECCo</i>	-9.33	7.34
vanilla	0.01	<i>Generic</i>	-8.89	6.88
vanilla	0.01	<i>Omniscient</i>	-8.67	6.87
vanilla	0.01	<i>REVISE</i>	-8.65	6.8
full	0.05	<i>ECCo</i>	-1.29	0.397
<b>full</b>	<b>0.05</b>	<b>Generic</b>	<b>-1.21</b>	<b>0.356</b>
full	0.05	<i>Omniscient</i>	-5.08	5.09
full	0.05	<i>REVISE</i>	$-5.91 \cdot 10^{27}$	$1.36 \cdot 10^{28}$
vanilla	0.05	<i>ECCo</i>	-9.35	7.32
vanilla	0.05	<i>Generic</i>	-8.85	6.87
vanilla	0.05	<i>Omniscient</i>	-8.7	6.96
vanilla	0.05	<i>REVISE</i>	-8.52	6.76
<b>full</b>	<b>0.1</b>	<b>ECCo</b>	<b>-1.2</b>	<b>0.383</b>
full	0.1	<i>Generic</i>	-1.5	0.735
full	0.1	<i>Omniscient</i>	-5.17	5.23
full	0.1	<i>REVISE</i>	$-3.06 \cdot 10^{26}$	$7.7 \cdot 10^{26}$
vanilla	0.1	<i>ECCo</i>	-9.33	7.32
vanilla	0.1	<i>Generic</i>	-8.88	6.86
vanilla	0.1	<i>Omniscient</i>	-8.69	6.9

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
vanilla	0.1	<i>REVISE</i>	-8.68	6.81
<b>full</b>	<b>0.5</b>	<b>ECCo</b>	<b>-1.12</b>	<b>0.217</b>
full	0.5	<i>Generic</i>	-1.21	0.352
full	0.5	<i>Omniscient</i>	-5.09	5.12
full	0.5	<i>REVISE</i>	$-5.97 \cdot 10^{27}$	$1.37 \cdot 10^{28}$
vanilla	0.5	<i>ECCo</i>	-9.35	7.3
vanilla	0.5	<i>Generic</i>	-8.89	6.92
vanilla	0.5	<i>Omniscient</i>	-8.68	6.93
vanilla	0.5	<i>REVISE</i>	-8.53	6.75
<b>full</b>	<b>1</b>	<b>ECCo</b>	<b>-1.1</b>	<b>0.163</b>
full	1	<i>Generic</i>	-1.49	0.726
full	1	<i>Omniscient</i>	-5.16	5.2
full	1	<i>REVISE</i>	$-3.09 \cdot 10^{26}$	$7.22 \cdot 10^{26}$
vanilla	1	<i>ECCo</i>	-9.34	7.36
vanilla	1	<i>Generic</i>	-8.86	6.85
vanilla	1	<i>Omniscient</i>	-8.7	6.9
vanilla	1	<i>REVISE</i>	-8.69	6.85
full	5	<i>ECCo</i>	-1.75	0.154
<b>full</b>	<b>5</b>	<b>Generic</b>	<b>-1.21</b>	<b>0.363</b>
full	5	<i>Omniscient</i>	-5.14	5.16
full	5	<i>REVISE</i>	$-1.1 \cdot 10^{28}$	$2.5 \cdot 10^{28}$
vanilla	5	<i>ECCo</i>	-9.36	7.32
vanilla	5	<i>Generic</i>	-8.88	6.91
vanilla	5	<i>Omniscient</i>	-8.7	6.93
vanilla	5	<i>REVISE</i>	-8.52	6.73
full	10	<i>ECCo</i>	$-1.02 \cdot 10^6$	$2.32 \cdot 10^6$
<b>full</b>	<b>10</b>	<b>Generic</b>	<b>-1.49</b>	<b>0.702</b>
full	10	<i>Omniscient</i>	-5.13	5.16
full	10	<i>REVISE</i>	$-3.74 \cdot 10^{26}$	$9.09 \cdot 10^{26}$
vanilla	10	<i>ECCo</i>	-9.31	7.33
vanilla	10	<i>Generic</i>	-8.87	6.86
vanilla	10	<i>Omniscient</i>	-8.7	6.89
vanilla	10	<i>REVISE</i>	-8.69	6.83
full	15	<i>ECCo</i>	$-3.31 \cdot 10^{13}$	$7.54 \cdot 10^{13}$
<b>full</b>	<b>15</b>	<b>Generic</b>	<b>-1.22</b>	<b>0.37</b>
full	15	<i>Omniscient</i>	-5.2	5.23
full	15	<i>REVISE</i>	$-9.01 \cdot 10^{27}$	$2.06 \cdot 10^{28}$
vanilla	15	<i>ECCo</i>	-9.38	7.34
vanilla	15	<i>Generic</i>	-8.86	6.87
vanilla	15	<i>Omniscient</i>	-8.69	6.96
vanilla	15	<i>REVISE</i>	-8.51	6.73