

000 001 002 003 004 005 DEEPWAVERL: SELF-SUPERVISED FULL WAVEFORM 006 INVERSION VIA REINFORCEMENT LEARNING 007 008 009

010 **Anonymous authors**
011 Paper under double-blind review
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029

ABSTRACT

030 Full Waveform Inversion (FWI) is a fundamental technique to estimate subsurface
031 geophysical properties, such as velocity, from seismic measurements. While su-
032 pervised deep learning methods have recently shown promising performance by
033 directly mapping seismic data to velocity maps, they require ground-truth velocity
034 maps, which are costly and impractical to obtain at scale. A recent self-supervised
035 approach (UPFWI) removes this dependency by leveraging a differentiable for-
036 ward operator to reconstruct seismic data from predictions. However, in some
037 practical settings, the forward operator can only be accessed as a black box (e.g.,
038 legacy or commercial). Moreover, for complex scenarios, the operator can even
039 be non-differentiable. In this paper, we address this limitation (i.e., the depen-
040 dency on derivatives of forward operators) by introducing reinforcement learning
041 (RL) into self-supervised FWI. Our method, named DeepWaveRL, reformulates
042 FWI as a policy learning problem, where the model generates velocity maps as
043 actions, and the forward operator is used only to compute rewards. This design
044 avoids backpropagation through the forward operator, thus eliminating the need
045 to compute its derivatives. Furthermore, we identify key strategies to stabilize re-
046inforcement learning in this challenging setting. In the absence of ground-truth
047 labels and differentiable forward operators, our method achieves competitive per-
048 formance compared to supervised counterparts. We believe our approach provides
049 a more flexible solution for the FWI research community.
050

1 INTRODUCTION

051 Subsurface imaging is essential for characterizing geological structures and geophysical properties
052 (e.g., velocity and impedance), with applications in energy exploration, carbon capture and seque-
053 stration, and earthquake early warning systems. A central technique in this domain is Full Waveform
054 Inversion (FWI), which estimates subsurface velocity maps from seismic measurements. Typically,
055 seismic data are acquired through seismic surveys, where an array of receivers records reflected and
056 refracted seismic waves. These waves are generated by controlled sources. Mathematically, for an
057 isotropic medium with constant density, the velocity map and seismic measurements are connected
058 by the acoustic wave equation:

$$059 \nabla^2 p(x, z, t) - \frac{1}{v(x, z)^2} \frac{\partial^2 p(x, z, t)}{\partial t^2} = s(x, z, t) , \quad (1)$$

060 where x denotes the horizontal offset, z the depth, $p(x, z, t)$ the pressure wavefield at spatial location
061 (x, z) and time t , $v(x, z)$ the wave propagation velocity at (x, z) , $s(x, z, t)$ the source term, and ∇^2
062 the Laplacian operator. In practice, seismic data are often collected at the surface (i.e., $p(x, z =$
063 $0, t)$). While FWI has the potential to produce high-resolution velocity maps, the inverse problem
064 itself is inherently non-linear and ill-posed. In addition, conventional physics-driven approaches face
065 additional challenges, as they require intensive computation due to repeated forward simulations per
066 sample and exhibit strong sensitivity to noise and initial conditions. These challenges have motivated
067 growing interest in data-driven deep learning methods.

068 A majority of data-driven methods (Wu & Lin, 2019; Zhang et al., 2019; Jin et al., 2024) adopt a
069 supervised learning paradigm and formulate FWI as an image-to-image translation task. As shown
070 in Figure 1, deep neural networks are trained to directly learn the mapping from seismic data to

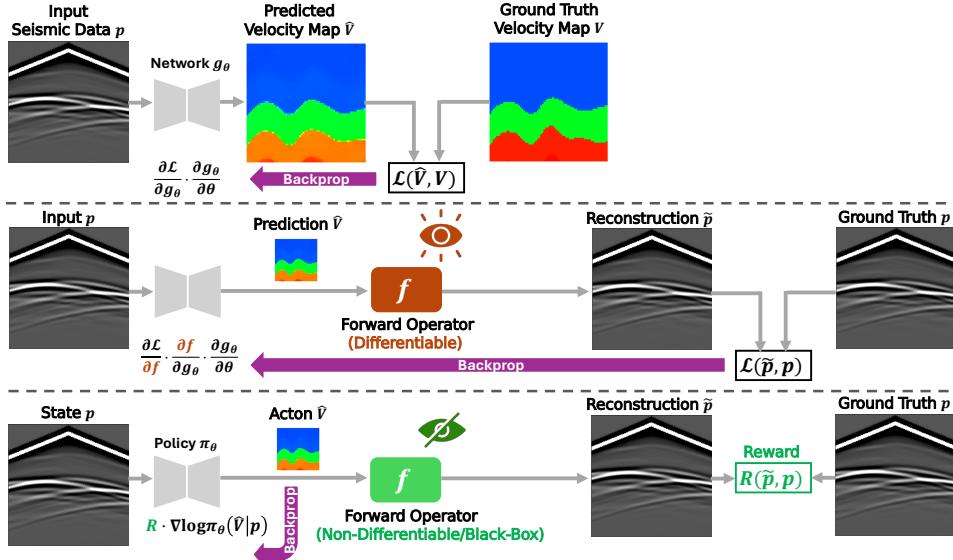


Figure 1: Comparison between different data-driven FWI methods. **Top:** Supervised learning methods compute loss between predicted and ground-truth velocity maps; **Middle:** Self-supervised learning method with a differentiable forward operator f computes the loss between input and reconstructed seismic data and backpropagates gradients through f ; **Bottom:** Our proposed DeepWaveRL uses the forward operator only to compute a reward signal based on misfit between seismic data, without backpropagating gradients through f , enabling greater flexibility.

velocity maps, enabling fast inference and achieving high accuracy under ideal conditions. However, these methods rely on a large amount of paired seismic data and velocity maps for training. In real-world scenarios, such ground-truth velocity maps are rarely available because constructing them is extremely time-consuming and requires substantial expertise from geophysicists.

A recent work (UPFWI, Jin et al., 2022) explicitly leverages the underlying physics knowledge and achieves self-supervised learning without ground-truth velocity maps. As illustrated in Figure 1, a *differentiable* forward operator f is coupled with a neural network to simulate seismic data from predicted velocity maps. By minimizing reconstruction loss on seismic data with *gradients backpropagated through f* , the network can be trained in an end-to-end manner without labeled supervision. However, this design imposes several critical limitations. First, the need for differentiability restricts the choice of forward solvers: many high-performance seismic simulators are implemented in low-level languages such as Fortran or C++ and are only available as non-differentiable “black boxes.” Second, real physical systems often exhibit non-smooth behaviors—for example, fractures that open only beyond a pressure threshold—where the wavefield response can change abruptly, violating differentiability and further limiting the applicability of such approaches.

In this paper, we present DeepWaveRL, a novel self-supervised approach for FWI that removes the dependency on differentiable forward operators by leveraging reinforcement learning (RL). As depicted in Figure 1, we formulate FWI as a single-step decision problem, where the input seismic data serve as the state, and a policy network outputs the corresponding velocity map as the action. A forward operator f is still used, but *only to compute rewards* based on the simulated seismic data from the predicted velocity map. The policy is then optimized via a policy gradient algorithm, which computes the gradients of the network’s action probabilities, weighted by the reward signal. Therefore, our method *eliminates the need to backpropagate gradients through f* .

Policy optimization in this setting poses unique challenges: continuous velocity values lead to an enormous action space, hindering effective exploration; large amplitude disparities between waves bias learning toward dominant ones; reward signals can only be evaluated for the entire velocity map without pixel-level feedback.

To address these issues, we further identify three key strategies. First, we use a discrete action space by partitioning the velocity range into finite bins, which significantly reduces the burden of exploration while maintaining accuracy. Second, we adopt a sign-preserving logarithmic transformation

108 for seismic data that compresses dominant directive wave energy and amplifies weaker signals (e.g.,
 109 reflections and deep arrivals), thereby yielding more precise predictions in deeper regions. Third, we
 110 exploit the ability of a well-trained policy to adapt across datasets, allowing transfer of knowledge
 111 in scenarios where training from scratch would be difficult or unstable.

112 We evaluate our method on several datasets from OpenFWI (Deng et al., 2022), a large-scale, multi-
 113 structural dataset collection. Experimental results show that our DeepWaveRL attains comparable
 114 performance to the supervised baseline InversionNet (Wu & Lin, 2019; Jin et al., 2024) on CurveVel-
 115 A, with a Mean Absolute Error (MAE) of 0.0527 (vs. 0.0409), a Root Mean Squared Error (RMSE)
 116 of 0.1012 (vs. 0.0944), and a Structured Similarity (SSIM) of 0.8601 (vs. 0.8796). DeepWaveRL
 117 with transfer learning also yields competitive performance on FlatFault-A and CurveFault-A.

118 Our contribution is summarized as follows:
 119

- 120 • We propose DeepWaveRL, a reinforcement learning framework for self-supervised full
 121 waveform inversion (FWI), which removes the need for differentiable forward operators.
- 122 • We propose three key techniques for stable and efficient training of DeepWaveRL, in-
 123 cluding discretized velocity actions, sign-preserving logarithmic transformation on seismic
 124 data, and transfer learning of well-trained policies.
- 125 • We demonstrate that our proposed DeepWaveRL achieves competitive performance with-
 126 out the involvement of ground-truth labels and differentiable forward operators.

128 2 METHOD

130 In this section, we first briefly summarize the state-of-the-art group-based reinforcement learning
 131 algorithms and then present our DeepWaveRL and its components. After that, we provide a summary
 132 of the comparison of DeepWaveRL with previous FWI methods from a gradient perspective.

134 2.1 PRELIMINARY

136 Shao et al. (2024) introduces Group Relative Policy Optimization (GRPO) that enhances Proximal
 137 Policy Optimization (PPO, Schulman et al., 2017) by omitting the value function and estimating
 138 the advantage in a group-relative manner. This is followed by several variants such as Decoupled
 139 Clip and Dynamic Sampling Policy Optimization (DAPO, Yu et al., 2025) and Group Sequence
 140 Policy Optimization (GSPO, Zheng et al., 2025), yielding even superior training efficiency and per-
 141 formance. The core idea of GRPO is summarized as follows.

142 For a specific question $q \sim P(Q)$, a group of G responses $\{o_i\}_{i=1}^G$ are sampled from an old policy
 143 network $\pi_{\theta_{\text{old}}}$. Each response o_i is then fed into a reward function to obtain the individual reward
 144 R_i . By normalizing the rewards within each group, an advantage \hat{A}_i is assigned to each response.
 145 The policy network is optimized by maximizing the following clipped objective, similar to PPO:

$$146 \mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)} \\ 147 \quad 148 \quad 149 \quad 150 \quad 151 \quad 152 \quad 153 \quad 154 \quad 155 \quad 156 \quad 157 \quad 158 \quad 159 \quad 160 \quad 161 \quad 162 \quad 163 \quad 164 \quad 165 \quad 166 \quad 167 \quad 168 \quad 169 \quad 170 \quad 171 \quad 172 \quad 173 \quad 174 \quad 175 \quad 176 \quad 177 \quad 178 \quad 179 \quad 180 \quad 181 \quad 182 \quad 183 \quad 184 \quad 185 \quad 186 \quad 187 \quad 188 \quad 189 \quad 190 \quad 191 \quad 192 \quad 193 \quad 194 \quad 195 \quad 196 \quad 197 \quad 198 \quad 199 \quad 200 \quad 201 \quad 202 \quad 203 \quad 204 \quad 205 \quad 206 \quad 207 \quad 208 \quad 209 \quad 210 \quad 211 \quad 212 \quad 213 \quad 214 \quad 215 \quad 216 \quad 217 \quad 218 \quad 219 \quad 220 \quad 221 \quad 222 \quad 223 \quad 224 \quad 225 \quad 226 \quad 227 \quad 228 \quad 229 \quad 230 \quad 231 \quad 232 \quad 233 \quad 234 \quad 235 \quad 236 \quad 237 \quad 238 \quad 239 \quad 240 \quad 241 \quad 242 \quad 243 \quad 244 \quad 245 \quad 246 \quad 247 \quad 248 \quad 249 \quad 250 \quad 251 \quad 252 \quad 253 \quad 254 \quad 255 \quad 256 \quad 257 \quad 258 \quad 259 \quad 260 \quad 261 \quad 262 \quad 263 \quad 264 \quad 265 \quad 266 \quad 267 \quad 268 \quad 269 \quad 270 \quad 271 \quad 272 \quad 273 \quad 274 \quad 275 \quad 276 \quad 277 \quad 278 \quad 279 \quad 280 \quad 281 \quad 282 \quad 283 \quad 284 \quad 285 \quad 286 \quad 287 \quad 288 \quad 289 \quad 290 \quad 291 \quad 292 \quad 293 \quad 294 \quad 295 \quad 296 \quad 297 \quad 298 \quad 299 \quad 300 \quad 301 \quad 302 \quad 303 \quad 304 \quad 305 \quad 306 \quad 307 \quad 308 \quad 309 \quad 310 \quad 311 \quad 312 \quad 313 \quad 314 \quad 315 \quad 316 \quad 317 \quad 318 \quad 319 \quad 320 \quad 321 \quad 322 \quad 323 \quad 324 \quad 325 \quad 326 \quad 327 \quad 328 \quad 329 \quad 330 \quad 331 \quad 332 \quad 333 \quad 334 \quad 335 \quad 336 \quad 337 \quad 338 \quad 339 \quad 340 \quad 341 \quad 342 \quad 343 \quad 344 \quad 345 \quad 346 \quad 347 \quad 348 \quad 349 \quad 350 \quad 351 \quad 352 \quad 353 \quad 354 \quad 355 \quad 356 \quad 357 \quad 358 \quad 359 \quad 360 \quad 361 \quad 362 \quad 363 \quad 364 \quad 365 \quad 366 \quad 367 \quad 368 \quad 369 \quad 370 \quad 371 \quad 372 \quad 373 \quad 374 \quad 375 \quad 376 \quad 377 \quad 378 \quad 379 \quad 380 \quad 381 \quad 382 \quad 383 \quad 384 \quad 385 \quad 386 \quad 387 \quad 388 \quad 389 \quad 390 \quad 391 \quad 392 \quad 393 \quad 394 \quad 395 \quad 396 \quad 397 \quad 398 \quad 399 \quad 400 \quad 401 \quad 402 \quad 403 \quad 404 \quad 405 \quad 406 \quad 407 \quad 408 \quad 409 \quad 410 \quad 411 \quad 412 \quad 413 \quad 414 \quad 415 \quad 416 \quad 417 \quad 418 \quad 419 \quad 420 \quad 421 \quad 422 \quad 423 \quad 424 \quad 425 \quad 426 \quad 427 \quad 428 \quad 429 \quad 430 \quad 431 \quad 432 \quad 433 \quad 434 \quad 435 \quad 436 \quad 437 \quad 438 \quad 439 \quad 440 \quad 441 \quad 442 \quad 443 \quad 444 \quad 445 \quad 446 \quad 447 \quad 448 \quad 449 \quad 450 \quad 451 \quad 452 \quad 453 \quad 454 \quad 455 \quad 456 \quad 457 \quad 458 \quad 459 \quad 460 \quad 461 \quad 462 \quad 463 \quad 464 \quad 465 \quad 466 \quad 467 \quad 468 \quad 469 \quad 470 \quad 471 \quad 472 \quad 473 \quad 474 \quad 475 \quad 476 \quad 477 \quad 478 \quad 479 \quad 480 \quad 481 \quad 482 \quad 483 \quad 484 \quad 485 \quad 486 \quad 487 \quad 488 \quad 489 \quad 490 \quad 491 \quad 492 \quad 493 \quad 494 \quad 495 \quad 496 \quad 497 \quad 498 \quad 499 \quad 500 \quad 501 \quad 502 \quad 503 \quad 504 \quad 505 \quad 506 \quad 507 \quad 508 \quad 509 \quad 510 \quad 511 \quad 512 \quad 513 \quad 514 \quad 515 \quad 516 \quad 517 \quad 518 \quad 519 \quad 520 \quad 521 \quad 522 \quad 523 \quad 524 \quad 525 \quad 526 \quad 527 \quad 528 \quad 529 \quad 530 \quad 531 \quad 532 \quad 533 \quad 534 \quad 535 \quad 536 \quad 537 \quad 538 \quad 539 \quad 540 \quad 541 \quad 542 \quad 543 \quad 544 \quad 545 \quad 546 \quad 547 \quad 548 \quad 549 \quad 550 \quad 551 \quad 552 \quad 553 \quad 554 \quad 555 \quad 556 \quad 557 \quad 558 \quad 559 \quad 560 \quad 561 \quad 562 \quad 563 \quad 564 \quad 565 \quad 566 \quad 567 \quad 568 \quad 569 \quad 570 \quad 571 \quad 572 \quad 573 \quad 574 \quad 575 \quad 576 \quad 577 \quad 578 \quad 579 \quad 580 \quad 581 \quad 582 \quad 583 \quad 584 \quad 585 \quad 586 \quad 587 \quad 588 \quad 589 \quad 590 \quad 591 \quad 592 \quad 593 \quad 594 \quad 595 \quad 596 \quad 597 \quad 598 \quad 599 \quad 600 \quad 601 \quad 602 \quad 603 \quad 604 \quad 605 \quad 606 \quad 607 \quad 608 \quad 609 \quad 610 \quad 611 \quad 612 \quad 613 \quad 614 \quad 615 \quad 616 \quad 617 \quad 618 \quad 619 \quad 620 \quad 621 \quad 622 \quad 623 \quad 624 \quad 625 \quad 626 \quad 627 \quad 628 \quad 629 \quad 630 \quad 631 \quad 632 \quad 633 \quad 634 \quad 635 \quad 636 \quad 637 \quad 638 \quad 639 \quad 640 \quad 641 \quad 642 \quad 643 \quad 644 \quad 645 \quad 646 \quad 647 \quad 648 \quad 649 \quad 650 \quad 651 \quad 652 \quad 653 \quad 654 \quad 655 \quad 656 \quad 657 \quad 658 \quad 659 \quad 660 \quad 661 \quad 662 \quad 663 \quad 664 \quad 665 \quad 666 \quad 667 \quad 668 \quad 669 \quad 670 \quad 671 \quad 672 \quad 673 \quad 674 \quad 675 \quad 676 \quad 677 \quad 678 \quad 679 \quad 680 \quad 681 \quad 682 \quad 683 \quad 684 \quad 685 \quad 686 \quad 687 \quad 688 \quad 689 \quad 690 \quad 691 \quad 692 \quad 693 \quad 694 \quad 695 \quad 696 \quad 697 \quad 698 \quad 699 \quad 700 \quad 701 \quad 702 \quad 703 \quad 704 \quad 705 \quad 706 \quad 707 \quad 708 \quad 709 \quad 710 \quad 711 \quad 712 \quad 713 \quad 714 \quad 715 \quad 716 \quad 717 \quad 718 \quad 719 \quad 720 \quad 721 \quad 722 \quad 723 \quad 724 \quad 725 \quad 726 \quad 727 \quad 728 \quad 729 \quad 730 \quad 731 \quad 732 \quad 733 \quad 734 \quad 735 \quad 736 \quad 737 \quad 738 \quad 739 \quad 740 \quad 741 \quad 742 \quad 743 \quad 744 \quad 745 \quad 746 \quad 747 \quad 748 \quad 749 \quad 750 \quad 751 \quad 752 \quad 753 \quad 754 \quad 755 \quad 756 \quad 757 \quad 758 \quad 759 \quad 760 \quad 761 \quad 762 \quad 763 \quad 764 \quad 765 \quad 766 \quad 767 \quad 768 \quad 769 \quad 770 \quad 771 \quad 772 \quad 773 \quad 774 \quad 775 \quad 776 \quad 777 \quad 778 \quad 779 \quad 780 \quad 781 \quad 782 \quad 783 \quad 784 \quad 785 \quad 786 \quad 787 \quad 788 \quad 789 \quad 790 \quad 791 \quad 792 \quad 793 \quad 794 \quad 795 \quad 796 \quad 797 \quad 798 \quad 799 \quad 800 \quad 801 \quad 802 \quad 803 \quad 804 \quad 805 \quad 806 \quad 807 \quad 808 \quad 809 \quad 810 \quad 811 \quad 812 \quad 813 \quad 814 \quad 815 \quad 816 \quad 817 \quad 818 \quad 819 \quad 820 \quad 821 \quad 822 \quad 823 \quad 824 \quad 825 \quad 826 \quad 827 \quad 828 \quad 829 \quad 830 \quad 831 \quad 832 \quad 833 \quad 834 \quad 835 \quad 836 \quad 837 \quad 838 \quad 839 \quad 840 \quad 841 \quad 842 \quad 843 \quad 844 \quad 845 \quad 846 \quad 847 \quad 848 \quad 849 \quad 850 \quad 851 \quad 852 \quad 853 \quad 854 \quad 855 \quad 856 \quad 857 \quad 858 \quad 859 \quad 860 \quad 861 \quad 862 \quad 863 \quad 864 \quad 865 \quad 866 \quad 867 \quad 868 \quad 869 \quad 870 \quad 871 \quad 872 \quad 873 \quad 874 \quad 875 \quad 876 \quad 877 \quad 878 \quad 879 \quad 880 \quad 881 \quad 882 \quad 883 \quad 884 \quad 885 \quad 886 \quad 887 \quad 888 \quad 889 \quad 890 \quad 891 \quad 892 \quad 893 \quad 894 \quad 895 \quad 896 \quad 897 \quad 898 \quad 899 \quad 900 \quad 901 \quad 902 \quad 903 \quad 904 \quad 905 \quad 906 \quad 907 \quad 908 \quad 909 \quad 910 \quad 911 \quad 912 \quad 913 \quad 914 \quad 915 \quad 916 \quad 917 \quad 918 \quad 919 \quad 920 \quad 921 \quad 922 \quad 923 \quad 924 \quad 925 \quad 926 \quad 927 \quad 928 \quad 929 \quad 930 \quad 931 \quad 932 \quad 933 \quad 934 \quad 935 \quad 936 \quad 937 \quad 938 \quad 939 \quad 940 \quad 941 \quad 942 \quad 943 \quad 944 \quad 945 \quad 946 \quad 947 \quad 948 \quad 949 \quad 950 \quad 951 \quad 952 \quad 953 \quad 954 \quad 955 \quad 956 \quad 957 \quad 958 \quad 959 \quad 960 \quad 961 \quad 962 \quad 963 \quad 964 \quad 965 \quad 966 \quad 967 \quad 968 \quad 969 \quad 970 \quad 971 \quad 972 \quad 973 \quad 974 \quad 975 \quad 976 \quad 977 \quad 978 \quad 979 \quad 980 \quad 981 \quad 982 \quad 983 \quad 984 \quad 985 \quad 986 \quad 987 \quad 988 \quad 989 \quad 990 \quad 991 \quad 992 \quad 993 \quad 994 \quad 995 \quad 996 \quad 997 \quad 998 \quad 999 \quad 1000 \quad 1001 \quad 1002 \quad 1003 \quad 1004 \quad 1005 \quad 1006 \quad 1007 \quad 1008 \quad 1009 \quad 1010 \quad 1011 \quad 1012 \quad 1013 \quad 1014 \quad 1015 \quad 1016 \quad 1017 \quad 1018 \quad 1019 \quad 1020 \quad 1021 \quad 1022 \quad 1023 \quad 1024 \quad 1025 \quad 1026 \quad 1027 \quad 1028 \quad 1029 \quad 1030 \quad 1031 \quad 1032 \quad 1033 \quad 1034 \quad 1035 \quad 1036 \quad 1037 \quad 1038 \quad 1039 \quad 1040 \quad 1041 \quad 1042 \quad 1043 \quad 1044 \quad 1045 \quad 1046 \quad 1047 \quad 1048 \quad 1049 \quad 1050 \quad 1051 \quad 1052 \quad 1053 \quad 1054 \quad 1055 \quad 1056 \quad 1057 \quad 1058 \quad 1059 \quad 1060 \quad 1061 \quad 1062 \quad 1063 \quad 1064 \quad 1065 \quad 1066 \quad 1067 \quad 1068 \quad 1069 \quad 1070 \quad 1071 \quad 1072 \quad 1073 \quad 1074 \quad 1075 \quad 1076 \quad 1077 \quad 1078 \quad 1079 \quad 1080 \quad 1081 \quad 1082 \quad 1083 \quad 1084 \quad 1085 \quad 1086 \quad 1087 \quad 1088 \quad 1089 \quad 1090 \quad 1091 \quad 1092 \quad 1093 \quad 1094 \quad 1095 \quad 1096 \quad 1097 \quad 1098 \quad 1099 \quad 1100 \quad 1101 \quad 1102 \quad 1103 \quad 1104 \quad 1105 \quad 1106 \quad 1107 \quad 1108 \quad 1109 \quad 1110 \quad 1111 \quad 1112 \quad 1113 \quad 1114 \quad 1115 \quad 1116 \quad 1117 \quad 1118 \quad 1119 \quad 1120 \quad 1121 \quad 1122 \quad 1123 \quad 1124 \quad 1125 \quad 1126 \quad 1127 \quad 1128 \quad 1129 \quad 1130 \quad 1131 \quad 1132 \quad 1133 \quad 1134 \quad 1135 \quad 1136 \quad 1137 \quad 1138 \quad 1139 \quad 1140 \quad 1141 \quad 1142 \quad 1143 \quad 1144 \quad 1145 \quad 1146 \quad 1147 \quad 1148 \quad 1149 \quad 1150 \quad 1151 \quad 1152 \quad 1153 \quad 1154 \quad 1155 \quad 1156 \quad 1157 \quad 1158 \quad 1159 \quad 1160 \quad 1161 \quad 1162 \quad 1163 \quad 1164 \quad 1165 \quad 1166 \quad 1167 \quad 1168 \quad 1169 \quad 1170 \quad 1171 \quad 1172 \quad 1173 \quad 1174 \quad 1175 \quad 1176 \quad 1177 \quad 1178 \quad 1179 \quad 1180 \quad 1181 \quad 1182 \quad 1183 \quad 1184 \quad 1185 \quad 1186 \quad 1187 \quad 1188 \quad 1189 \quad 1190 \quad 1191 \quad 1192 \quad 1193 \quad 1194 \quad 1195 \quad 1196 \quad 1197 \quad 1198 \quad 1199 \quad 1200 \quad 1201 \quad 1202 \quad 1203 \quad 1204 \quad 1205 \quad 1206 \quad 1207 \quad 1208 \quad 1209 \quad 1210 \quad 1211 \quad 1212 \quad 1213 \quad 1214 \quad 1215 \quad 1216 \quad 1217 \quad 1218 \quad 1219 \quad 1220 \quad 1221 \quad 1222 \quad 1223 \quad 1224 \quad 1225 \quad 1226 \quad 1227 \quad 1228 \quad 1229 \quad 1230 \quad 1231 \quad 1232 \quad 1233 \quad 1234 \quad 1235 \quad 1236 \quad 1237 \quad 1238 \quad 1239 \quad 1240 \quad 1241 \quad 1242 \quad 1243 \quad 1244 \quad 1245 \quad 1246 \quad 1247 \quad 1248 \quad 1249 \quad 1250 \quad 1251 \quad 1252 \quad 1253 \quad 1254 \quad 1255 \quad 1256 \quad 1257 \quad 1258 \quad 1259 \quad 1260 \quad 1261 \quad 1262 \quad 1263 \quad 1264 \quad 1265 \quad 1266 \quad 1267 \quad 1268 \quad 1269 \quad 1270 \quad 1271 \quad 1272 \quad 1273 \quad 1274 \quad 1275 \quad 1276 \quad 1277 \quad 1278 \quad 1279 \quad 1280 \quad 1281 \quad 1282 \quad 1283 \quad 1284 \quad 1285 \quad 1286 \quad 1287 \quad 1288 \quad 1289 \quad 1290 \quad 1291 \quad 1292 \quad 1293 \quad 1294 \quad 1295 \quad 1296 \quad 1297 \quad 1298 \quad 1299 \quad 1300 \quad 1301 \quad 1302 \quad 1303 \quad 1304 \quad 1305 \quad 1306 \quad 1307 \quad 1308 \quad 1309 \quad 1310 \quad 1311 \quad 1312 \quad 1313 \quad 1314 \quad 1315 \quad 1316 \quad 1317 \quad 1318 \quad 1319 \quad 1320 \quad 1321 \quad 1322 \quad 1323 \quad 1324 \quad 1325 \quad 1326 \quad 1327 \quad 1328 \quad 1329 \quad 1330 \quad 1331 \quad 1332 \quad 1333 \quad 1334 \quad 1335 \quad 1336 \quad 1337 \quad 1338 \quad 1339 \quad 1340 \quad 1341 \quad 1342 \quad 1343 \quad 1344 \quad 1345 \quad 1346 \quad 1347 \quad 1348 \quad 1349 \quad 1350 \quad 1351 \quad 1352 \quad 1353 \quad 1354 \quad 1355 \quad 1356 \quad 1357 \quad 1358 \quad 1359 \quad 1360 \quad 1361 \quad 1362 \quad 1363 \quad 1364 \quad 1365 \quad 1366 \quad 1367 \quad 1368 \quad 1369 \quad 1370 \quad 1371 \quad 1372 \quad 1373 \quad 1374 \quad 1375 \quad 1376 \quad 1377 \quad 1378 \quad 1379 \quad 1380 \quad 1381 \quad 1382 \quad 1383 \quad 1384 \quad 1385 \quad 1386 \quad 1387 \quad 1388 \quad 1389 \quad 1390 \quad 1391 \quad 1392 \quad 1393 \quad 1394 \quad 1395 \quad 1396 \quad 1397 \quad 1398 \quad 1399 \quad 1400 \quad 1401 \quad 1402 \quad 1403 \quad 1404 \quad 1405 \quad 1406 \quad 1407 \quad 1408 \quad 1409 \quad 1410 \quad 1411 \quad 1412 \quad 1413 \quad 1414 \quad 1415 \quad 1416 \quad 1417 \quad 1418 \quad 1419 \quad 1420 \quad 1421 \quad 1422 \quad 1423 \quad 1424 \quad 1425 \quad 1426 \quad 1427 \quad 1428 \quad 1429 \quad 1430 \quad 1431 \quad 1432 \quad 1433 \quad 1434 \quad 1435 \quad 1436 \quad 1437 \quad 1438 \quad 1439 \quad 1440 \quad 1441 \quad 1442 \quad 1443 \quad 1444 \quad 1445 \quad 1446 \quad 1447 \quad 1448 \quad 1449 \quad 1450 \quad 1451 \quad 1452 \quad 1453 \quad 1454 \quad 1455 \quad 1456 \quad 1457 \quad 1458 \quad 1459 \quad 1460 \quad 1461 \quad 1462 \quad 1463 \quad 1464 \quad 1465 \quad 1466 \quad 1467 \quad 1468 \quad 1469 \quad 1470 \quad 1471 \quad 1472 \quad 1473 \quad 1474 \quad 1475 \quad 1476 \quad 1477 \quad 1478 \quad 1479 \quad 1480 \quad 1481 \quad 1482 \quad 1483 \quad 1484 \quad 1485 \quad 1486 \quad 1487 \quad 1488 \quad 1489 \quad 1490 \quad 1491 \quad 1492 \quad 1493 \quad 1494 \quad 1495 \quad 1496 \quad 1497 \quad 1498 \quad 1499 \quad 1500 \quad 1501 \quad 1502 \quad 1503 \quad 1504 \quad 1505 \quad 1506 \quad 1507 \quad 1508 \quad 1509 \quad 1510 \quad 1511 \quad 1512 \quad 1513 \quad 1514 \quad 1515 \quad 1516 \quad 1517 \quad 1518 \quad 1519 \quad 1520 \quad 1521 \quad 1522 \quad 1523 \quad 1524 \quad 1525 \quad 1526 \quad 1527 \quad 1528 \quad 1529 \quad 1530 \quad 1531 \quad 1532 \quad 1533 \quad 1534 \quad 1535 \quad 1536 \quad 1537 \quad 1538 \quad 1539 \quad 1540 \quad 1541 \quad 1542 \quad 1543 \quad 1544 \quad 1545 \quad 1546 \quad 1547 \quad 1548 \quad 1549 \quad 1550 \quad 1551 \quad 1552 \quad 1553 \quad 1554 \quad 1555 \quad 1556 \quad 1557 \quad 1558 \quad 1559 \quad 1560 \quad 1561 \quad 1562 \quad 1563 \quad 1564 \quad 1565 \quad 1566 \quad 1567 \quad 1568 \quad 1569 \quad 1570 \quad 1571 \quad 1572 \quad 1573 \quad 1574 \quad 1575 \quad 1576 \quad 1577 \quad 1578 \quad 1579 \quad 1580 \quad 1581 \quad 1582 \quad 1583 \quad 1584 \quad 1585 \quad 1586 \quad 1587 \quad 1588 \quad 1589 \quad 1590 \quad 1591 \quad 1592 \quad 1593 \quad 1594 \quad 1595 \quad 1596 \quad 1597 \quad 1598 \quad 1599 \quad 1600 \quad 1601 \quad 1602 \quad 1603 \quad 1604 \quad 1605 \quad 1606 \quad 1607 \quad 1608 \quad 1609 \quad 1610 \quad 1611 \quad 1612 \quad 1613 \quad 1614 \quad 1615 \quad 1616 \quad 1617 \quad 1618 \quad 1619 \quad 1620 \quad 1621 \quad 1622 \quad 1623 \quad 1624 \quad 1625 \quad 1626 \quad 1627 \quad 1628 \quad 1629 \quad 1630 \quad 1631 \quad 1632 \quad 1633 \quad 1634 \quad 1635 \quad 1636 \quad 1637 \quad 1638 \quad 1639 \quad 1640 \quad 1641 \quad 1642 \quad 1643 \quad 1644 \quad 1645 \quad 1646 \quad 1647 \quad 1648 \quad 1649 \quad 1650 \quad 1651 \quad 1652 \quad 1653 \quad 1654 \quad 1655 \quad 1656 \quad 1657 \quad 1658 \quad 1659 \quad 1660 \quad 1661 \quad 1662 \quad 1663 \quad 1664 \quad 1665 \quad 1666 \quad 1667 \quad 1668 \quad 1669 \quad 1670 \quad 1671 \quad 1672 \quad 1673 \quad 1674 \quad 1675 \quad 1676 \quad 1677 \quad 1678 \quad 1679 \quad 1680 \quad 1681 \quad 1682 \quad 1683 \quad 1684 \quad 1685 \quad 1686 \quad 1687 \quad 1688 \quad 1689 \quad 1690 \quad 1691 \quad 1692 \quad 1693 \quad 1694 \quad 1695 \quad 1696 \quad 1697 \quad 1698 \quad 1699 \quad 1700 \quad 1701 \quad 1$$

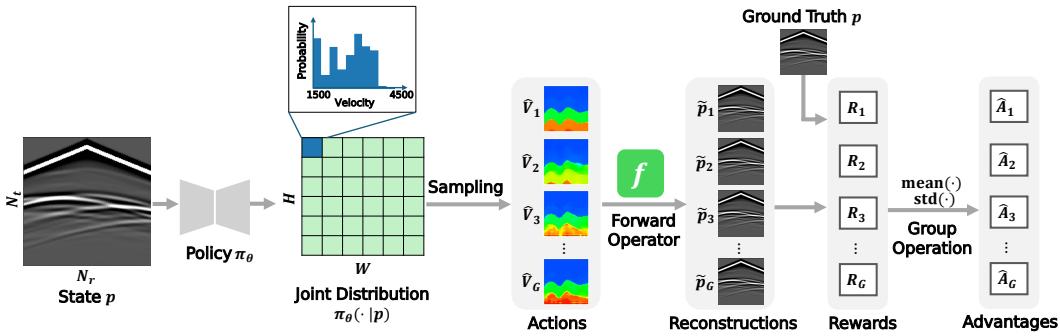


Figure 2: Schematic illustration of the policy optimization pipeline in DeepWaveRL. For brevity, we only show the seismic data from one source.

Here, N_s denotes the number of sources used during data acquisition, N_t the number of recorded timesteps, and N_r the total number of receivers. The network then outputs the probability distribution of velocity $\pi_\theta(\cdot | p)_{h,w}$ at each spatial location (h, w) . Following GRPO, we sample a group of G velocity maps $\{\hat{V}_i\}_{i=1}^G$ from the predicted joint probability distribution as 2D actions, where each $\hat{V}_i = \{\hat{v}_{i,h,w}\}_{h=1}^H, w=1}^W$, and H and W are the vertical and horizontal dimensions of the velocity map. To reduce exploration burden, we use a discrete action space instead of the continuous one.

To assign a reward R_i to each sampled action \hat{V}_i , we employ a forward operator f to simulate seismic data $\tilde{p}_i = f(\hat{V}_i) \in \mathbb{R}^{N_s \times N_t \times N_r}$. The reward is then computed based on the misfit between the input seismic data p and the reconstruction \tilde{p}_i . We further compute the relative advantage \hat{A}_i within each group following Equation 3, but use a map-level importance ratio.

According to a recent work (Zheng et al., 2025), the mismatch between the unit of reward and the unit of optimization objective can introduce high-variance noise and further lead to model collapse. In our settings, a reward is assigned to the whole 2D velocity map. Therefore, instead of computing the pixel-level importance ratios in Equation 3, we define the map-level importance ratio as:

$$m_i = \left[\frac{\pi_\theta(\hat{V}_i | p)}{\pi_{\theta_{\text{old}}}(\hat{V}_i | p)} \right]^{\frac{1}{H \cdot W}} = \exp \left[\frac{1}{H \cdot W} \sum_{h=1}^H \sum_{w=1}^W \log \left(\frac{\pi_\theta(\hat{v}_{i,h,w} | p)}{\pi_{\theta_{\text{old}}}(\hat{v}_{i,h,w} | p)} \right) \right]. \quad (4)$$

Consequently, the overall optimization objective can now be written as:

$$\begin{aligned} \mathcal{J}(\theta) = & \mathbb{E}_{p \sim \mathcal{P}, \{\hat{V}_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | p)} \\ & \frac{1}{G} \sum_{i=1}^G \left\{ \min \left[m_i(\theta) \hat{A}_i, \text{clip} \left(m_i(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}} \right) \hat{A}_i \right] \right\}, \end{aligned} \quad (5)$$

where we follow Yu et al. (2025) and decouple the lower and higher clipping range as ε_{low} and $\varepsilon_{\text{high}}$, and \mathcal{P} denotes the distribution of seismic data. We also remove the KL penalty term as our initial model is not as good as common pretrained language models, and we allow the model distribution to diverge from the initial model.

The policy network is thus trained to shift its output distribution toward higher-reward actions. Importantly, the entire training process is self-supervised, without the involvement of ground-truth velocity maps. Furthermore, no gradients are backpropagated through the forward operator f , which allows f to be arbitrary, including non-differentiable or black-box simulators.

Additionally, we note that the policy network can also be optimized at test time, since only seismic data are required. This test-time optimization further boosts performance, which will be discussed in Section 3.3.

2.3 KEY STRATEGIES FOR STABLE AND EFFICIENT LEARNING

Discrete Action Space. During training, we treat the predicted velocity map $\hat{V} = \{\hat{v}_{h,w}\}_{h=1}^H, w=1}^W$ as a 2D action sampled from the joint distribution. We initially experimented with a continuous action

space, where the network outputs two maps of size $H \times W$, representing the mean $\mu_{h,w}$ and standard deviation $\sigma_{h,w}$ of a Gaussian distribution at each spatial location. Hence, the velocity at location (h, w) was sampled as $\hat{v}_{h,w} \sim \mathcal{N}(\mu_{h,w}, \sigma_{h,w})$. However, this formulation led to unstable training and noisy predictions. To address these issues, we discretize the velocity values uniformly into B bins, treating the sampled action as a predicted category. Given an action $a_{h,w} \in [0, 1, \dots, B - 1]$, we compute the corresponding velocity as

$$v(a_{h,w}) = \frac{v_{\max} - v_{\min}}{B} \cdot (a_{h,w} + 0.5) + v_{\min}, \quad (6)$$

where v_{\max} and v_{\min} are the largest and smallest possible velocities in a dataset. This discrete action space substantially improves both training stability and prediction quality.

Sign-Preserving Logarithm Transformation in Reward. Similar to the loss function of UP-FWI (Jin et al., 2022), we define the reward as the negative pixel-wise ℓ_1 and ℓ_2 distance between the input and reconstructed seismic data. To further enhance learning, we apply a sign-preserving logarithm transformation during the computation of rewards as:

$$p' = \text{sign}(p) \cdot \log(k \cdot |p| + c), \quad (7)$$

where k and c are hyperparameters to control the strength of the transformation. This non-linear transformation can compress dominant directive wave energy and amplify weaker signals (e.g., reflections and deep arrivals), thereby guiding the network to recover more accurate velocities in deeper regions. The reward function can then be described as:

$$R = -\ell_1(p', \tilde{p}') - \ell_2(p', \tilde{p}'). \quad (8)$$

Transfer Learning using Well-Trained Policies. Directly training our DeepWaveRL from scratch sometimes leads to unstable learning and convergence issues. We provide examples of predicted velocity maps generated by these models in Appendix D.5. The predictions exhibit unrealistically low velocities in deep regions after certain training steps. To address this issue, we propose to initialize the policy network with weights from a well-trained model on a different dataset. This transfer learning strategy allows prior knowledge of velocity distributions to be reused across datasets, leading to more stable training and improved convergence.

2.4 COMPARISON WITH PREVIOUS FWI METHODS FROM A GRADIENT PERSPECTIVE

To demonstrate the relationship among supervised, self-supervised with a differentiable forward operator, and RL-based self-supervised approaches, we provide analysis from the perspective of gradient construction and propagation.

Supervised: Gradients are directly computed from discrepancies with ground-truth velocity maps, $\nabla_{\theta} \mathcal{L}_{\text{sup}} \propto \frac{\partial \|V - \hat{V}\|}{\partial \hat{V}}$. While this yields stable optimization and strong supervision, it is impractical in real-world scenarios due to the scarcity of paired data.

Differentiable self-supervised: Gradients originate from seismic reconstruction error and back-propagate through a differentiable forward operator, $\nabla_{\theta} \mathcal{L}_{\text{diff}} \propto \frac{\partial f(\hat{V})}{\partial \hat{V}}$, enabling physics-informed learning but restricted by differentiability and high computational cost.

RL-based self-supervised: In contrast to self-supervised methods with a differentiable forward operator, where the loss must be differentiable with respect to \hat{V} and thus optimization is tightly coupled to f , our RL-based method replaces this requirement by converting non-differentiable errors into reward signals that act as multipliers of the policy gradients, $\nabla_{\theta} \mathcal{J}(\theta) \propto \mathbb{E} [\hat{A} \cdot \nabla_{\theta} \log m_i(\theta)]$.

Full derivations and detailed comparisons are provided in Appendix B.

3 EXPERIMENTS

In this section, we evaluate the performance of our proposed DeepWaveRL on OpenFWI (Deng et al., 2022), comparing it with both supervised and self-supervised baselines. We also investigate the impact of the discrete action space and examine the effect of the logarithm transformation on seismic data through ablation studies.

270
271
272

Table 1: Quantitative results evaluated on CurveVel-A. For models with discrete predictions, we report the mean estimate.

273
274
275
276
277
278

Method	MAE \downarrow	RMSE \downarrow	SSIM \uparrow
InversionNet	0.0409	0.0944	0.8796
UPFWI	0.0805	0.1411	0.8443
DeepWaveRL	0.0717	0.1300	0.8303
DeepWaveRL + TTO	<u>0.0527</u>	<u>0.1012</u>	<u>0.8601</u>

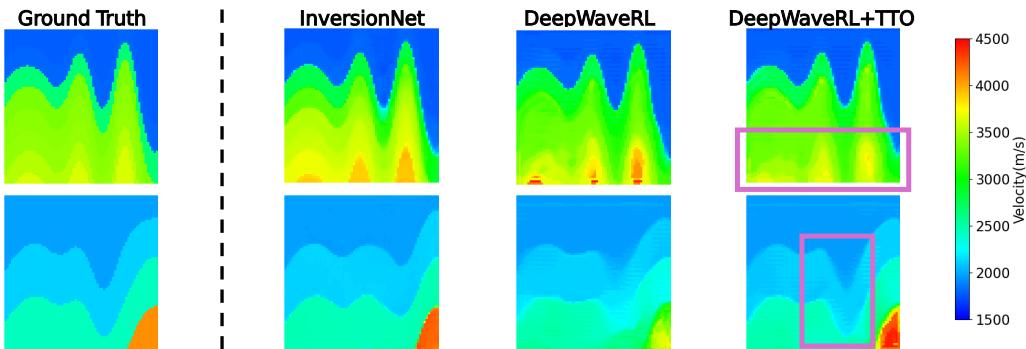
279
280
281
282
283
284
285
286
287
288
289
290
291

Figure 3: Illustration of ground truth and inversion results of different methods on CurveVel-A.

293

294 3.1 DATASETS

295
296
297
298
299
300
301
302
303

We verify our method on CurveVel-A, FlatFault-A and CurveFault-A of OpenFWI (Deng et al., 2022), an open-source collection of large-scale, multi-structural benchmark datasets for data-driven seismic FWI. CurveVel-A contains velocity maps composed of curved layers with clear interfaces, while FlatFault-A and CurveFault-A focus more on geological fault identification and have velocity maps with flat and curved layers, respectively. CurveVel-A contains 30K velocity maps and their corresponding seismic data. Following the official data split, we use 24K for training and 6K for testing. FlatFault-A and CurveFault-A contain 54K samples each, and we use 48K/6K splitting.

304
305
306
307
308
309
310
311

Each velocity map in all three datasets has a size of 70×70 , with a grid size of 10 meters in both horizontal and depth directions. The velocity value ranges from 1,500 meter/second to 4,500 meter/second. For seismic data, five sources are placed evenly with a 170-meter spacing and a central source frequency of 15 Hz. The seismic data are recorded by 70 receivers at 10-meter intervals, each collecting 1,000 timesteps over 1 second. This results in seismic data of shape $5 \times 1000 \times 70$. For additional details, we refer readers to the original OpenFWI paper (Deng et al., 2022).

312
313
314
315
316
317
318
319

Evaluation Metrics: We evaluate predicted velocity maps using MAE, RMSE, and Structural Similarity (SSIM), consistent with prior work (Wu & Lin, 2019; Feng et al., 2024; Deng et al., 2022). MAE and RMSE quantify pixel-wise errors, while SSIM captures perceptual similarity, reflecting the structured information of velocity maps where distortions can be easily perceived by a human. Note that all measurements are computed on normalized velocity maps, with MAE and RMSE in the range $[-1, 1]$, and SSIM in $[0, 1]$.

320
321
322
323

Comparison: We compare our method to InversionNet (Wu & Lin, 2019) which achieves the state-of-the-art performance when trained solely on each dataset, as demonstrated in a recent work (Jin et al., 2024). Additionally, we list the benchmarking results of UPFWI (Jin et al., 2022) from the original OpenFWI paper.

Technical details regarding training are provided in Appendix C.

324

325

Table 2: Quantitative results evaluated on FlatFault-A and CurveFault-A.

326

327

328

329

330

331

332

333

334

335

336

Dataset	Method	MAE \downarrow	RMSE \downarrow	SSIM \uparrow
FlatFault-A	InversionNet	0.0098	0.0276	0.9880
	UPFWI	0.0876	0.2060	0.9340
	DeepWaveRL	0.0301	0.0557	0.9062
	DeepWaveRL + TTO	0.0268	0.0476	0.9146
CurveFault-A	InversionNet	0.0164	0.0480	0.9721
	UPFWI	0.0500	0.0966	0.9495
	DeepWaveRL	0.0362	0.0703	0.9111
	DeepWaveRL + TTO	0.0278	0.0502	0.9323

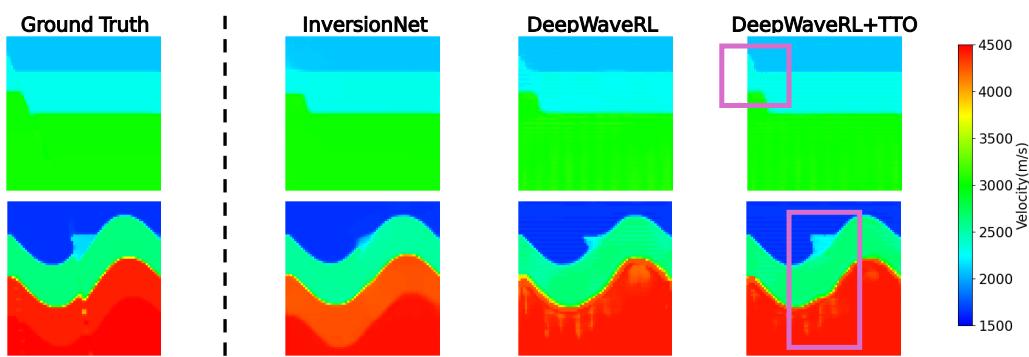


Figure 4: Illustration of ground truth and inversion results of different methods on FlatFault-A (top) and CurveFault-A (bottom).

351

352

353

3.3 MAIN RESULTS

354

Results on CurveVel-A: Table 1 shows the quantitative results of different methods on CurveVel-A. For the models that predict velocity as discrete values, we report the mean prediction (expected value) where we compute the expectation over the bin values using their predicted probabilities.

Among all the models, InversionNet yields the best performance, which is expected as it leverages supervised learning and directly predicts continuous velocity values. In comparison, our DeepWaveRL with test-time optimization (DeepWaveRL+TTO) attains the second-best performance with a slight gap. Notably, test-time optimization substantially boosts the performance of DeepWaveRL.

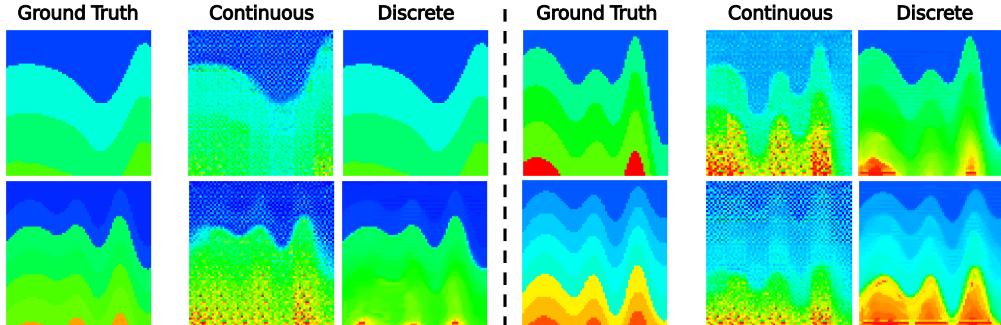
Figure 3 further illustrates examples of ground-truth velocity maps and inversion results from different methods. Here, the results of DeepWaveRL models are all mean predictions. Consistent with our quantitative analysis, InversionNet produces sharp layer boundaries and smooth and uniform velocities within layers. However, in certain regions, DeepWaveRL+TTO yields more details. For instance, as highlighted in the first row, the predictions of DeepWaveRL+TTO have more accurate velocities in deep regions, whereas InversionNet introduces additional layers and predicts inaccurate velocities. Another observation is that only DeepWaveRL+TTO precisely reconstructs the subsurface structure in the highlighted areas in the second row. Moreover, when comparing the results of DeepWaveRL and DeepWaveRL+TTO, we find that test-time optimization helps eliminate artifacts, recover curved structures, and improve the accuracy of intra-layer velocities. More visualization results are shown in Appendix D.1.

Transfer Learning Results on FlatFault-A and CurveFault-A: Table 2 lists the quantitative results on FlatFault-A and CurveFault-A, and mean predictions are reported. For transfer learning, we adopt the last checkpoint of the DeepWaveRL model trained with test-time optimization on CurveVel-A. For both datasets, our DeepWaveRL+TTO consistently outperforms UPFWI in terms of MAE and RMSE, with small gaps in SSIM. This indicates that DeepWaveRL-TTO yields generally accurate predictions, but there may be small shifts or artifacts that are visually noticeable.

378
379
380
381
382
383

Table 3: Quantitative results evaluated on CurveVel-A, with different choices of action space.

Setting	MAE \downarrow	RMSE \downarrow	SSIM \uparrow
Continuous Action Space	0.1443	0.2073	0.4359
Discrete Action Space	0.0527	0.1012	0.8601

395
396 Figure 5: Illustration of ground truth and inversion results of DeepWaveRL with different choices
397 of action space398
399
400 The visualization results in Figure 4 further support our hypothesis. Despite some stripe artifacts,
401 our DeepWaveRL and DeepWaveRL+TTO generate more precise details in some shallow regions.
402 As highlighted in the first row (FlatFault-A), DeepWaveRL models reconstruct the fault on the top-
403 left corner, while it is barely visible in the predictions of InversionNet. Similarly, in the second
404 row (CurveFault-A), our DeepWaveRL models precisely capture the triangle-shaped region. More
405 visualization results are shown in Appendix D.2.406 3.4 ABLATIONS
407408 **Continuous vs. Discrete Action Space:** We further analyse how the choice of action space affects
409 model performance. The quantitative results are summarized in Table 3. For simplicity, we denote
410 DeepWaveRL with a discrete action space as DeepWaveRL-D, and with a continuous action space
411 as DeepWaveRL-C. In terms of all three metrics, DeepWaveRL-D outperforms DeepWaveRL-C to a
412 large extent. The gap is particularly large in SSIM, suggesting that DeepWaveRL-C fails to produce
413 results consistent with human perceptual quality.414 Figure 5 provides qualitative comparisons between DeepWave-C and DeepWaveRL-D, which fur-
415 ther support our quantitative analysis. While DeepWaveRL-C can recover some structures in shallow
416 regions, there are plenty of artifacts all over the predictions. In particular, the high-velocity areas
417 are severely corrupted, making boundaries unrecognizable. Furthermore, these artifacts persist even
418 when training is extended, indicating that the continuous action space poses significant optimiza-
419 tion challenges. By contrast, discretization reduces the complexity of the action space and greatly
420 stabilizes training, leading to more accurate and reliable results.421 **With or Without Logarithm Transformation:** To evaluate the effect of the sign-preserving log-
422 arithm transformation, we train our DeepWaveRL without transformation on CurveVel-A. The re-
423 sulting MAE, RMSE, and SSIM are 0.0785, 0.1383, and 0.8195, respectively. Compared to the
424 results of DeepWaveRL with transformation in Table 1, the performance degrades in terms of all
425 three metrics. This is consistent with the visualization results in Appendix D.3, where some details
426 are missing in deep regions.427 4 DISCUSSION
428430 During qualitative analysis, we find that mode collapse occurs in some of the predictions of our
431 DeepWaveRL, as illustrated in Appendix D.4. The ground-truth velocity maps of these predictions
have close velocities in their shallow regions, but this pattern does not guarantee the occurrence of

432 mode collapse. Thus, we may take this into consideration in our future work. Another limitation
 433 is that our DeepWaveRL still struggles with the recovery of the structures in deeper regions due to
 434 the inherent attenuation of signals in these regions. Our sign-preserving logarithm transformation is
 435 one of the solutions, but more advanced algorithms are still needed. Furthermore, our DeepWaveRL
 436 framework enables another potential research direction, which is to incorporate non-differentiable
 437 regularization terms such as total variation into the reward design.

440 5 RELATED WORK

441
 442 **Deep Learning for FWI:** Deep learning approaches to FWI span data-driven, physics-informed,
 443 and hybrid paradigms (Lin et al., 2023; Adler et al., 2021; Yu & Ma, 2021). Fully supervised meth-
 444 ods (Araya-Polo et al., 2018; Wu & Lin, 2019; Zhang et al., 2019; Li et al., 2020) learn direct
 445 mappings from seismic data to velocity models using paired data, which are costly to acquire and
 446 often lead to poor generalization under domain shifts. To reduce reliance on labels and improve
 447 robustness, self-supervised strategies have emerged. Feng et al. (Feng et al., 2022; 2024) decouple
 448 the seismic encoder and velocity decoder by leveraging latent space correlations, enabling separate
 449 training. SiameseFWI (Saad et al., 2024) explores self-supervision with a Siamese network that bet-
 450 ter aligns simulated and observed data. Semi-supervised learning has also been explored. Sun et al.
 451 (2023) proposes a CycleGAN-based framework to reconstruct missing low-frequency components
 452 in field data. Other methods generate pseudo-labels from unlabeled or auxiliary data (Rojas-Gómez
 453 et al., 2022; Cai et al., 2022), bridging the gap between labeled and unlabeled domains. Unsuper-
 454 vised methods such as UPFWI (Jin et al., 2022) and Jia et al. (Jia et al., 2025) go further by eliminat-
 455 ing labels entirely. These approaches minimize waveform mismatches under physical constraints,
 456 using differentiable forward modeling to optimize predicted velocity maps. However, their reliance
 457 on computationally intensive and differentiable solvers limits scalability to high-resolution, elastic,
 458 or 3D FWI, and precludes use with black-box simulators. Diffusion models offer an alternative by
 459 learning generative priors that guide the inverse process via plug-and-play (PnP) denoising (Song
 460 et al., 2022; Chung et al., 2023; Zhang et al., 2025). Wang et al. (2023) successfully applies this strat-
 461 egy to FWI. While these models avoid paired supervision, they still require large velocity datasets
 462 for training and incur high inference costs due to repeated forward simulations during sampling.

463 **Reinforcement learning and group-based policy optimization:** Policy-gradient RL methods (e.g.,
 464 PPO (Schulman et al., 2017)) provide a principled way to optimize stochastic policies via likelihood-
 465 ratio estimators, and have been widely used in sequential generation and control. Recent advances in
 466 group-/sequence-level policy optimization demonstrate that performing importance-weighting and
 467 clipping at the unit-of-reward level (group or sequence) can reduce variance and stabilize training for
 468 structured outputs. In particular, DeepSeekMath (Shao et al., 2024) replaces value-function estima-
 469 tion with group-relative normalization, yielding more efficient and stable updates. They show that
 470 this framework can be scaled to mathematical reasoning tasks with strong generalization. DAPO (Yu
 471 et al., 2025) extends the paradigm by decoupling clipping ranges and introducing dynamic sampling,
 472 further improving stability under diverse reward distributions. GSPO (Zheng et al., 2025) general-
 473 izes these ideas to full sequence-level optimization with the importance ratio based on sequence
 474 likelihood, aligning long-horizon objectives with token-level policies. Inspired by this progression,
 475 DeepWaveRL adapts the same philosophy to the geophysics domain by treating an entire velocity
 476 map as a structured action, making map-level optimization a natural extension of group-based RL
 477 methods for physics-driven inverse problems.

478 6 CONCLUSION

479
 480 In this study, we introduce DeepWaveRL, a reinforcement learning framework for self-supervised
 481 full waveform inversion that eliminates the need for differentiable forward operators. By incorpo-
 482 rating discretized velocity actions, a sign-preserving logarithmic transformation of seismic data, and
 483 transfer learning from well-trained policies, DeepWaveRL achieves stable and efficient training. We
 484 demonstrate through experiments that our method attains competitive performance without relying
 485 on ground-truth velocity maps or differentiable forward operators. This approach provides a flexible
 486 solution for FWI, offering new possibilities in real-world settings.

486 REFERENCES
487

488 Amir Adler, Mauricio Araya-Polo, and Tomaso Poggio. Deep learning for seismic inverse problems:
489 Toward the acceleration of geophysical analysis workflows. *IEEE Signal Processing Magazine*,
490 38:89–119, 2021. doi: <https://doi.org/10.1109/MSP.2020.3037429>.

491 Mauricio Araya-Polo, Joseph Jennings, Amir Adler, and Taylor Dahlke. Deep-learning tomography.
492 *The Leading Edge*, 37(1):58–66, 2018.

493

494 Ao Cai, Hongrui Qiu, and Fenglin Niu. Semi-supervised surface wave tomography with wasser-
495 stein cycle-consistent GAN: method and application to southern california plate boundary re-
496 gion. *Journal of Geophysical Research: Solid Earth*, 127(3):e2021JB023598, 2022. doi:
497 <https://doi.org/10.1029/2021JB023598>.

498 Hyungjin Chung, Jeongsol Kim, Michael T McCann, Marc L Klasky, and Jong Chul Ye. Diffu-
499 sion posterior sampling for general noisy inverse problems. In *11th International Conference on*
500 *Learning Representations, ICLR 2023*, 2023.

501

502 Chengyuan Deng, Shihang Feng, Hanchen Wang, Xitong Zhang, Peng Jin, Yinan Feng, Qili Zeng,
503 Yinpeng Chen, and Youzuo Lin. Openfwi: Large-scale multi-structural benchmark datasets for
504 full waveform inversion. *Advances in Neural Information Processing Systems*, 35:6007–6020,
505 2022.

506 Yinan Feng, Yinpeng Chen, Shihang Feng, Peng Jin, Zicheng Liu, and Youzuo Lin. An intriguing
507 property of geophysics inversion. In *International Conference on Machine Learning*, pp. 6434–
508 6446. PMLR, 2022.

509

510 Yinan Feng, Yinpeng Chen, Peng Jin, Shihang Feng, and Youzuo Lin. Auto-linear phenomenon in
511 subsurface imaging. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org, 2024.

512

513 Anqi Jia, Jian Sun, Bo Du, and Yuzhao Lin. Seismic full waveform inversion with uncertainty
514 analysis using unsupervised variational deep learning. *IEEE Transactions on Geoscience and*
515 *Remote Sensing*, 63:1–16, 2025. doi: 10.1109/TGRS.2025.3564647.

516

517 Peng Jin, Xitong Zhang, Yinpeng Chen, Sharon Xiaolei Huang, Zicheng Liu, and Youzuo Lin. Unsu-
518 pervised learning of full-waveform inversion: Connecting CNN and partial differential equation in
519 a loop. In *Proceedings of the Tenth International Conference on Learning Representations (ICLR)*,
520 2022.

521

522 Peng Jin, Yinan Feng, Shihang Feng, Hanchen Wang, Yinpeng Chen, Benjamin Consolvo, Zicheng
523 Liu, and Youzuo Lin. An empirical study of large-scale data-driven full waveform inversion.
524 *Scientific Reports*, 14(1):20034, 2024.

525

526 Shucai Li, Bin Liu, Yuxiao Ren, Yangkang Chen, Senlin Yang, Yunhai Wang, and Peng Jiang. Deep-
527 Learning Inversion of Seismic Data. *IEEE Transactions on Geoscience and Remote Sensing*, 58
528 (3):2135–2149, March 2020. ISSN 1558-0644. doi: 10.1109/TGRS.2019.2953473.

529

530 Youzuo Lin, James Theiler, and Brendt Wohlberg. Physics-guided data-driven seismic inversion:
531 Recent progress and future opportunities in full waveform inversion. *IEEE Signal Processing*
532 *Magazine*, 40:115–133, 2023. doi: <https://doi.org/10.1109/MSP.2022.3217658>.

533

534 Renán Rojas-Gómez, Jihyun Yang, Youzuo Lin, James Theiler, and Brendt Wohlberg. Physics-
535 consistent data-driven waveform inversion with adaptive data augmentation. *IEEE Geoscience*
536 *and Remote Sensing Letters*, 19:1–5, 2022. doi: <https://doi.org/10.1109/LGRS.2020.3022021>.

537

538 Omar M. Saad, Randy Harsuko, and Tariq Alkhalifah. SiameseFWI: a deep learning network for
539 enhanced full waveform inversion. *Journal of Geophysical Research: Machine Learning and*
540 *Computation*, 1(3):e2024JH000227, 2024. doi: <https://doi.org/10.1029/2024JH000227>.

541

542 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
543 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

540 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
 541 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical
 542 reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.

543

544 Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging
 545 with score-based generative models. In *International Conference on Learning Representations*,
 546 2022.

547 Hongyu Sun, Yen Sun, Rami Nammourand, Christian Rivera, Paul Williamson, and Laurent De-
 548 manet. Learning with real data without real labels: a strategy for extrapolated full-waveform
 549 inversion with field data. *Geophysical Journal International*, 235:1761–1777, 2023. doi:
 550 <https://doi.org/10.1093/gji/ggad330>.

551

552 Fu Wang, Xinquan Huang, and Tariq Alkhalifah. A prior regularized full waveform inversion using
 553 generative diffusion models. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–11,
 554 2023. doi: <https://doi.org/10.1109/TGRS.2023.3337014>.

555

556 Yue Wu and Youzuo Lin. InversionNet: An efficient and accurate data-driven full waveform inver-
 557 sion. *IEEE Transactions on Computational Imaging*, 6:419–433, 2019.

558

559 Qiyi Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian
 560 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system
 561 at scale. *arXiv preprint arXiv:2503.14476*, 2025.

562

563 Xiwei Yu and Jianwei Ma. Deep learning for geophysics: Current and future trends. *Reviews of
 564 Geophysics*, 59:e2021RG000742, 2021. doi: <https://doi.org/10.1029/2021RG000742>.

565

566 Bingliang Zhang, Wenda Chu, Julius Berner, Chenlin Meng, Anima Anandkumar, and Yang Song.
 567 Improving diffusion inverse problem solving with decoupled noise annealing. In *Proceedings of
 568 the Computer Vision and Pattern Recognition Conference*, pp. 20895–20905, 2025.

569

570 Zhongping Zhang, Yue Wu, Zheng Zhou, and Youzuo Lin. VelocityGAN: Subsurface velocity
 571 image estimation using conditional adversarial networks. In *2019 IEEE Winter Conference on
 572 Applications of Computer Vision (WACV)*, pp. 705–714. IEEE, 2019.

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594 APPENDIX
595596 A USE OF LARGE LANGUAGE MODELS (LLMs)
597598
599 In preparing this manuscript, we utilized a large language model (LLM) to assist with the writing
600 process. Its role was limited to improving language quality, including grammar, phrasing, and over-
601 all readability, as well as helping with L^AT_EX formatting for tables and equations. The conception
602 of the research, design of experiments, analysis of results, and all scientific contributions were the
603 responsibility of the authors.
604605 B COMPARISON WITH PREVIOUS FWI METHODS FROM A GRADIENT
606 PERSPECTIVE
607608 To elucidate the relationship among supervised, self-supervised with a differentiable forward op-
609 erator, and reinforcement learning based self-supervised approaches, we provide a unified analysis
610 from the perspective of gradient construction and propagation.
611612 **Supervised learning:** In the supervised paradigm, a neural network g_θ maps seismic measure-
613 ments p to a predicted velocity map $\hat{V} = g_\theta(p)$. Training relies on paired ground-truth velocity
614 maps v , with a loss function of the form
615

616
$$\mathcal{L}_{\text{sup}}(\theta) = \mathbb{E}_{(p, V) \sim \mathcal{D}} \|V - g_\theta(p)\|, \quad (9)$$

617 where \mathcal{D} denotes the joint distribution of seismic data and corresponding velocity maps, and $\|\cdot\|$
618 denotes a generic norm (e.g., ℓ_1 , ℓ_2 , or mixed/perceptual norms). The gradient is obtained via direct
619 backpropagation:
620

621
$$\nabla_\theta \mathcal{L}_{\text{sup}} = \mathbb{E}_{(p, V) \sim \mathcal{D}} \frac{\partial \mathcal{L}_{\text{sup}}}{\partial \hat{V}} \cdot \frac{\partial \hat{V}}{\partial \theta} = \mathbb{E}_{(p, V) \sim \mathcal{D}} \underbrace{\frac{\partial \mathcal{L}_{\text{sup}}}{\partial g_\theta(p)}}_{\text{Ground truth involved}} \cdot \frac{\partial g_\theta(p)}{\partial \theta}, \quad (10)$$

625 where the learning signal is explicitly anchored to the availability of ground-truth velocity maps.
626 While this yields stable optimization and strong supervision, it is impractical in real-world scenarios
627 due to the scarcity of paired data.
628629 **Self-supervised learning with a differentiable forward operator:** The UPFWI framework (Jin
630 et al., 2022) removes the dependence on ground truth by incorporating a differentiable forward
631 operator f that simulates seismic data from predicted velocity maps:
632

633
$$\hat{V} = g_\theta(p), \quad \tilde{p} = f(\hat{V}). \quad (11)$$

634 The reconstruction objective is defined as
635

636
$$\mathcal{L}_{\text{diff}}(\theta) = \mathbb{E}_{p \sim \mathcal{P}} \|p - \tilde{p}\| = \mathbb{E}_{p \sim \mathcal{P}} \|p - f(g_\theta(p))\|, \quad (12)$$

638 with gradients computed via the chain rule:
639

640
$$\nabla_\theta \mathcal{L}_{\text{diff}} = \mathbb{E}_{p \sim \mathcal{P}} \frac{\partial \mathcal{L}_{\text{diff}}}{\partial \tilde{p}} \cdot \frac{\partial \tilde{p}}{\partial \hat{V}} \cdot \frac{\partial \hat{V}}{\partial \theta} = \mathbb{E}_{p \sim \mathcal{P}} \frac{\partial \mathcal{L}_{\text{diff}}}{\partial f(g_\theta(p))} \cdot \underbrace{\frac{\partial f(g_\theta(p))}{\partial g_\theta(p)}}_{\text{Differentiable}} \cdot \frac{\partial g_\theta(p)}{\partial \theta}. \quad (13)$$

643 This formulation enables end-to-end training using only seismic data, yet a differentiable forward
644 operator is required to compute $\partial f / \partial \hat{V}$.
645646 **Self-supervised learning via reinforcement learning:** Our proposed DeepWaveRL relaxes the
647 differentiability constraint by reframing FWI as a policy optimization problem. With the definition

648 in the above sections, we can derive the gradient of our objective as follows (clipping is omitted for
 649 brevity):
 650

$$651 \nabla_{\theta} \mathcal{J}(\theta) = \nabla_{\theta} \mathbb{E}_{p \sim \mathcal{P}, \{\hat{V}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | p)} \frac{1}{G} \sum_{i=1}^G \left\{ m_i(\theta) \hat{A}_i \right\} \quad (14)$$

$$654 = \mathbb{E}_{p \sim \mathcal{P}, \{\hat{V}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | p)} \left[\frac{1}{G} \sum_{i=1}^G m_i(\theta) \hat{A}_i \cdot \nabla_{\theta} \log m_i(\theta) \right] \quad (15)$$

$$656 = \mathbb{E}_{p \sim \mathcal{P}, \{\hat{V}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | p)} \quad (16)$$

$$658 \left[\frac{1}{G} \sum_{i=1}^G \left(\frac{\pi_{\theta}(\hat{V}_i | p)}{\pi_{\theta_{\text{old}}}(\hat{V}_i | p)} \right)^{\frac{1}{H \cdot W}} \underbrace{\hat{A}_i}_{\text{Forward model only involved as a multiplier}} \cdot \frac{1}{H \cdot W} \sum_{h,w} \nabla_{\theta} \log \pi_{\theta}(\hat{v}_{i,h,w} | p) \right]. \quad (17)$$

663 Unlike supervised and differentiable self-supervised learning, where the error must be differentiable
 664 with respect to \hat{V} and thus couples optimization tightly to the properties of f , reinforcement learning
 665 replaces this requirement by transforming non-differentiable errors into reward signals that directly
 666 reweight policy gradients.
 667

668 Thus, the three paradigms can be interpreted within a common gradient-based framework: **Supervised**:
 669 Gradients are directly computed from discrepancies with labeled velocity maps, yielding
 670 stable optimization but requiring costly ground truth. **Differentiable self-supervised**: Gradients
 671 originate from seismic reconstruction error and backpropagate through a differentiable forward op-
 672 erator, enabling physics-informed learning but restricted by differentiability and high computational
 673 cost. **RL-based self-supervised**: Gradients arise from log-likelihood weighting in policy space,
 674 with seismic misfit entering only as a reward. This bypasses differentiability, accommodates arbi-
 675 trary forward operators, and enables stochastic exploration in complex inversion landscapes.
 676

C TECHNICAL DETAILS

678 We normalize the input seismic data to the range $[-1, 1]$ and apply the logarithm transformation
 679 with $k = 3$ and $c = 0$ on seismic data. For optimization, we employ the AdamW optimizer with
 680 momentum parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a weight decay of 1×10^{-4} to update all
 681 parameters of the network. The details of hyperparameters and training settings are provided in
 682 Table 4. The ϵ_{low} and ϵ_{high} are 0.2 and 0.27, respectively. For the network architecture, we adopt
 683 a four-layer encoder-decoder Vision Transformer (ViT), and we append four convolutional blocks
 684 with upsampling layers ($5 \times$ and $2 \times$), batch normalization, and leaky ReLU as activation functions
 685 to map the output of the decoder to 70×70 velocity map with 100 bins. We implement our models
 686 in Pytorch and train them on 16 NVIDIA H100 GPUs.
 687

	CVA		FFA & CFA	
Test-time Optimization		✓		✓
Training Steps	44,880	1,440	7,360	1,600
Initial Learning Rate	8e-4	1.6e-4	6.4e-4	6.4e-4
Learning Rate Decay	/	1,360	/	/
Batch Size	128	2048	2048	2048
Group Size	256	16	32	32

696 Table 4: Training details
 697
 698

D VISUALIZATIONS

D.1 MORE VISUALIZATION RESULTS ON CURVEVEL-A

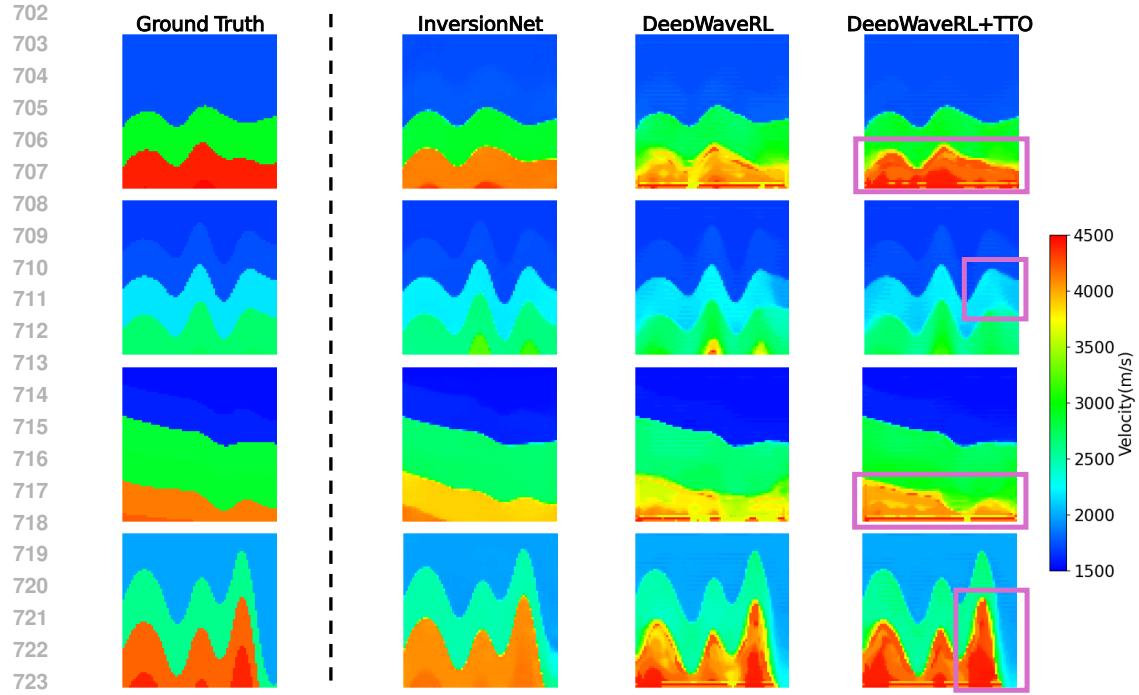


Figure 6: Illustration of ground truth and inversion results of different methods on CurveVel-A.

D.2 MORE VISUALIZATION RESULTS ON FLATFAULT-A AND CURVEFAULT-A

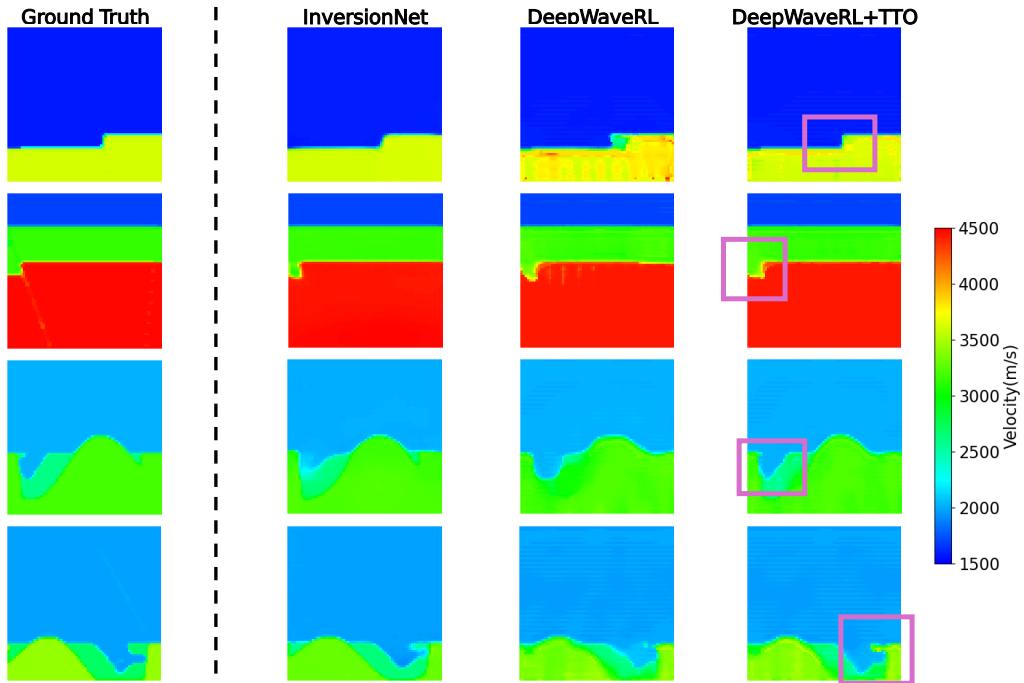
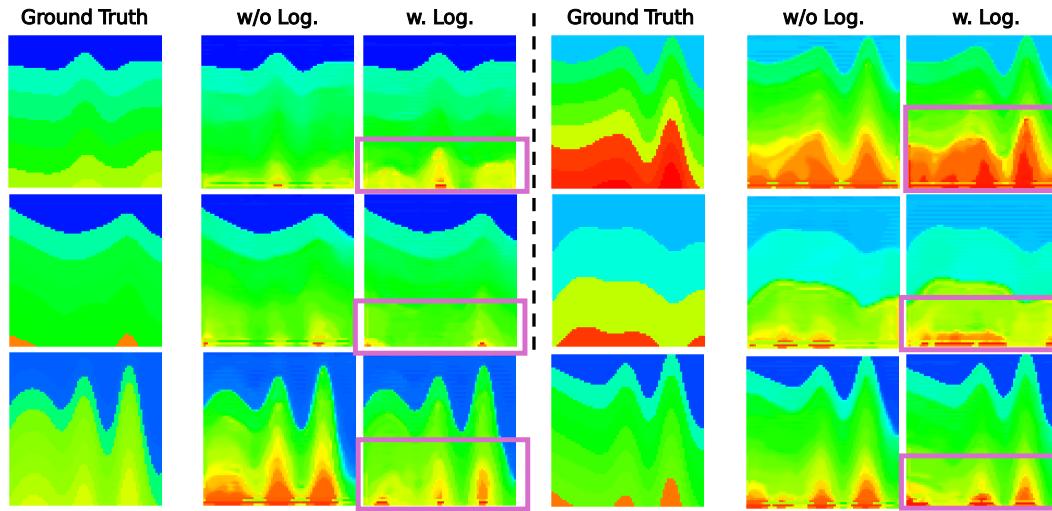
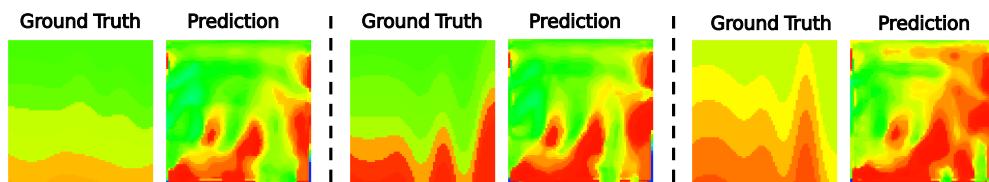
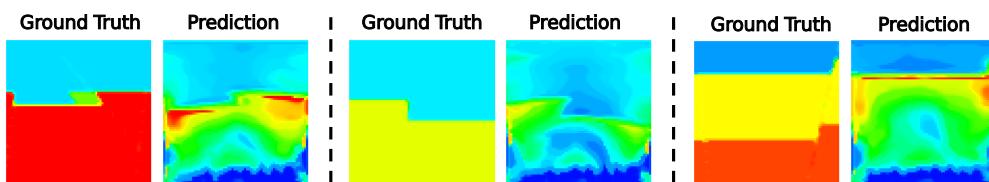


Figure 7: Illustration of ground truth and inversion results of different methods on FlatFault-A and CurveFault-A.

756 D.3 VISUALIZATION RESULTS OF DEEPWAVERL WITHOUT LOGARITHM TRANSFORMATION
757758
759 Figure 8: Illustration of ground truth and inversion results of DeepWaveRL with and without the
760 sign-preserving logarithm transformation
761
762763 D.4 FAILURE CASES
764765
766 Figure 9: Examples of failure cases where the predictions collapse to similar patterns.
767
768769 D.5 EXAMPLES OF UNSTABLE TRAINING RESULTS
770771
772 Figure 10: Examples of predicted velocity maps generated by the model that has experienced unsta-
773 ble training.
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809