
Simplifying Full Waveform Inversion via Domain-Independent Self-Supervised Learning

Yinan Feng

Earth and Environmental Sciences Division
Los Alamos National Laboratory
Los Alamos, NM 87544, USA
ynf@lanl.gov

Yinpeng Chen

Microsoft Corporation
Redmond WA 98052, USA
yiche@microsoft.com

Peng Jin

College of Information Sciences and Technology
The Pennsylvania State University
University Park, PA, 16802, USA
pqj5125@psu.edu

Shihang Feng

Earth and Environmental Sciences Division
Los Alamos National Laboratory
Los Alamos, NM 87544, USA
shihang@lanl.gov

Zicheng Liu

Microsoft Corporation
Redmond WA 98052, USA
zliu@microsoft.com

Youzuo Lin

Earth and Environmental Sciences Division
Los Alamos National Laboratory
Los Alamos, NM 87544, USA
ylin@lanl.gov

Abstract

Geophysics has witnessed success in applying deep learning to one of its core problems: full waveform inversion (FWI) to predict subsurface velocity maps from seismic data. It is treated as an image-to-image translation problem, jointly training an encoder for seismic data and a decoder for the velocity map from seismic-velocity pairs. In this paper, we report a surprising phenomenon: when training an encoder and decoder separately in their own domains via self-supervised learning, a linear relationship is observed across domains in the latent spaces. Moreover, this phenomenon connects multiple FWI datasets in an elegant manner: these datasets can share the self-learned encoder and decoder with different linear mappings.

Based on these findings, we develop SimFWI, a new paradigm that includes two-steps: (a) learning a seismic encoder and a velocity decoder separately by masked image modeling over multiple datasets; (b) learning a linear mapping per dataset. Experimental results show that SimFWI can achieve comparable results to a jointly trained model from the supervision of paired seismic data and velocity maps.

1 Introduction

Geophysical inversion techniques are crucial for revealing subsurface layering and geophysical properties (such as velocity and conductivity), supporting important applications such as energy exploration, carbon capture, and sequestration, groundwater contamination and remediation, and earthquake early warning systems. In this field, the full waveform inversion (FWI) is a well-known method to infer subsurface velocity map from the seismic data, which are mathematically connected

by an acoustic wave equation as:

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

where $p(x, z, t)$ represents the seismic data and $c(x, z)$ is the velocity map. $s(x, z, t)$ is the source term. x is the horizontal direction, z is the depth, t denotes time, and ∇^2 is the Laplacian operator. In practice, seismic data is usually collected by the sensors on the surface. Thus, we only have the 2D seismic data $p(x, z = 0, t)$, abbreviating it as $p(x, t)$.

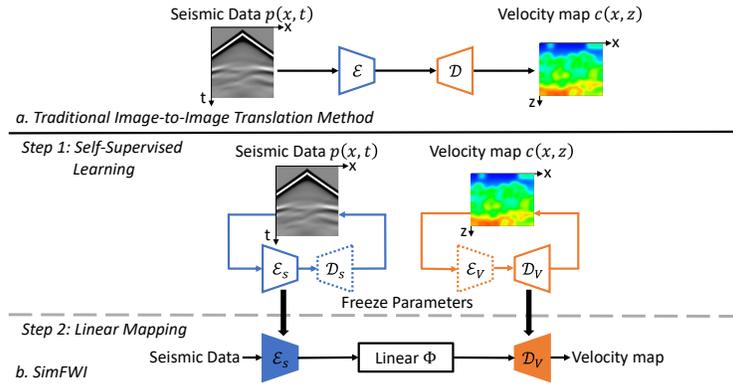


Figure 1: **Overview of SimFWI.** Compared to the jointly trained encoder-decoder (top), SimFWI (bottom) decouples the encoder and decoder and self-supervised trains them separately in their own domains. Then, a linear converter is learned to connect the frozen, pre-trained encoder and decoder.

Recent works Wu & Lin (2019); Zhang et al. (2019); Sun et al. (2021); Jin et al. (2022) consider FWI as an *image-to-image translation problem* constrained by a wave equation, and leverage deep neural networks in the solution. As shown in Figure 1-a, they learn an encoder-decoder architecture to map seismic data to velocity. Note that the encoder and decoder are jointly trained from the supervision of paired seismic data and velocity maps.

In this paper, we shift the paradigm to decouple the encoder and decoder and train them separately in their own domains via self-supervised learning. In particular, we train two masked autoencoders (MAE) He et al. (2022) separately, i.e. one for seismic data and one for velocity maps (see Figure 1-b). Surprisingly, we observe a linear correlation between the two latent spaces. This means the self-pretrained encoder and decoder can be frozen, and we only need to learn a linear converter to connect them from the paired seismic data and velocity map. This introduces an interesting insight into FWI: ***the self-consistent representation within each domain is associated with simpler mapping across domains.*** We name this method *SimFWI*, as it simplifies the mapping (linear) in FWI between seismic data and velocity map via domain-independent self-supervised learning.

Furthermore, SimFWI provides a better understanding of the relationship among multiple FWI datasets with different subsurface structures. We found that these datasets can share both encoders and decoders, but have different linear mappings between the latent spaces of two domains (i.e. seismic data and velocity map). Essentially, the two domains have a piece-wise linear relationship over multiple datasets. In addition, we found a correlation between the linear layer’s singular values and the complexity of the dataset.

SimFWI achieves solid performance on multiple FWI datasets. It has comparable results to the InversionNet Wu & Lin (2019), a jointly trained model that uses paired data as supervision, with only half the model size (12.3M vs. 24.4M). In the few-shot context where only limited paired data exists, SimFWI outperforms the InversionNet. Moreover, it is more robust to large, noisy data and the pre-trained encoder and decoder have a strong generalization ability.

Variable	Definition
$p(x, t)$	seismic data
$c(x, z)$	velocity maps
\mathcal{E}_s	encoder of seismic data
\mathcal{D}_s	decoder of seismic data
\mathcal{E}_v	encoder of velocity map
\mathcal{D}_v	decoder of velocity map
Φ	linear converter

Table 1: Table of Notation.

2 Related Works

Recently, data-driven methods for FWI have been developed. They consider the FWI as an image-to-image problem and jointly train the encoder-decoder network to solve it. Araya-Polo et al. (2018) use a fully connected network to invert velocity maps. Wu & Lin (2019) adopted an encoder-decoder CNN to solve. Zhang et al. (2019) employ GAN and transfer learning to improve the generalization. In Zeng et al. (2021), authors present an efficient and scalable encoder-decoder neural network for 3D FWI. Feng et al. (2021) develop a multi-scale framework with two convolutional neural networks to reconstruct the low- and high-frequency components of velocity maps. A thorough review of deep learning for FWI can be found in Lin et al. (2023).

Jin et al. (2022) use the finite difference to approximate the forward modeling as a differentiable operator and integrate it and a deep neural network (DNN) in a loop to construct an unsupervised learning method. Chen et al. (2021) proposed a self-supervised approach to solve the inverse problem from the perspective of image invariance. These purely self-supervised and unsupervised methods focus on how to solve problems without labels and still treat the network as a black box. Unlike them, our method uses self-supervised learning as a tool with the aim of simplifying the problem and decoupling the inverse process. We hope this can help the field better understand the problem and the relationship among different subsurface structures.

Recently, OpenFWI was released. It is the first open-source collection of large-scale multi-structural benchmark datasets for FWI Deng et al. (2022). It includes 12 datasets (11 2D datasets and one 3D dataset) synthesized from multiple sources. The datasets cover diverse domains in geophysics, such as interfaces, faults, and CO2 reservoirs, and feature a variety of subsurface structures, including flat and curved geologies. Along with the dataset, they also report performance benchmarks by using state-of-art data-driven methods and the physics-driven method.

3 Methodology

Recent works Wu & Lin (2019); Zhang et al. (2019); Sun et al. (2021); Jin et al. (2022) treat the full waveform inversion (FWI) as an image-to-image translation (from seismic data to velocity map) problem and leverage encoder-decoder architecture to achieve a significant performance boost. Here, we present new insights from the perspective of self-supervised learning. In particular, the encoder and decoder can be learned separately in their own domains via MAE He et al. (2022), and the two corresponding latent spaces are linearly correlated. Table 1 lists notations of our method.

3.1 Domain-Independent Self-Supervised Learning

We decouple seismic data $p(x, t)$ and velocity maps $c(x, z)$ and train individual masked autoencoders (MAE) He et al. (2022) for each domain (shown in Figure 1-b). Let us denote the encoder and decoder for seismic data as \mathcal{E}_s and \mathcal{D}_s . They are trained by using seismic data alone. Similarly, the encoder and decoder for the velocity map are denoted as \mathcal{E}_v and \mathcal{D}_v , which are trained by using velocity maps alone. As the pre-training is separate for the two domains, the pairing of seismic data and velocity maps is not needed at this stage. Note that the decoder for seismic data \mathcal{D}_s and the encoder for velocity map \mathcal{E}_v are auxiliary and will be removed. The encoder for seismic data \mathcal{E}_s and the decoder for the velocity map \mathcal{D}_v are the outputs of this domain-independent pre-training.

3.2 Two Intriguing Properties

We observe two intriguing properties in the latent spaces of two pre-trained models (seismic data and velocity maps). Note that the pre-trained models in both domains are *frozen*.

Property 1: The latent representation of two domains are linearly correlated. Surprisingly, without any fine-tuning, a simple relationship is observed between the latent embeddings of paired seismic data and velocity maps. This demonstrates an interesting association between self-consistent representation within each domain and simpler linear mapping across two domains.

Property 2: The encoder and decoder can be shared across datasets, while the linear mapping is dataset-specific. When dealing with multiple datasets with different subsurface structures, we find that the self-learned encoder and decoder can be shared across datasets (i.e. performing MAE on the combination of multiple datasets). The linear mapping (as mentioned in property 1) still holds for each dataset, but it is dataset-specific. Essentially, the two domains have a piece-wise linear relation globally over multiple datasets. Each dataset shows linear relation locally.

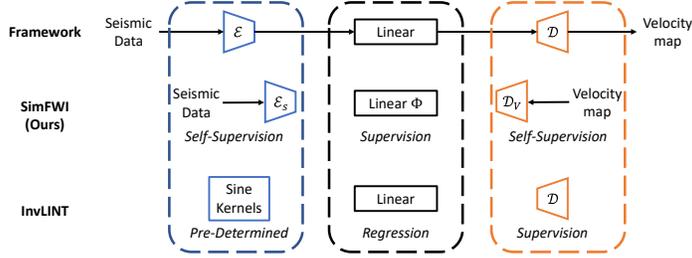


Figure 2: Comparison between our SimFWI and InvLINT Feng et al. (2022). The first row shows the similar framework of both methods. The second row indicates how our method trains each component. The third row shows how InvLINT trains each component.

3.3 SimFWI

As shown in Figure 1-b, SimFWI contains three components: a seismic encoder \mathcal{E}_s , a linear converter Φ , and a velocity decoder \mathcal{D}_v . They are trained in two steps: a domain-independent self-supervised learning step to train two autoencoders, and a supervised learning step to train the linear converter.

In the self-supervised learning step, as described above, we use Masked Autoencoder (MAE) He et al. (2022) as the self-supervised learners. In this step, paired data are not needed. Two MAEs are trained in their own domains. The trained \mathcal{E}_s and \mathcal{D}_v are frozen. Due to a lack of constraint, the self-supervised learner can easily learn shortcuts for reconstruction. We chose to use MAE because it generates better latent representations and can learn essential information about two physical quantities through the use of masks as noise, making it easier to connect the latent spaces of two modalities and transform and reconstruct the other.

In the supervised learning step, the converter Φ is trained to connect the frozen \mathcal{E}_s and \mathcal{D}_v with paired data. In practice, we decompose the linear converter into two linear layers with a low-dimensional bottleneck to constrain its rank. This effectively reduces redundancy in the network.

3.4 Comparing with the recent work

Compared with joint training methods (e.g., InversionNet), SimFWI decouple the training of the encoder and decoder. Only a linear converter is trained in a supervised manner with paired data. We achieve comparable results with only half the model size. Moreover, the pre-trained encoder and decoder have strong generalization ability and good handling of few-shot situations and noisy data.

In Feng et al. (2022), the authors also decouple the encoder and decoder, and use a linear layer to connect two latent spaces. They provided theoretical proof that establishes a near-linear relationship in a high-dimensional latent space when an appropriate encoder and decoder are used. But the difference is they use two pre-determined integral transforms, with Sine and Gaussian kernels, to embed the seismic data and velocity into high-dimensional spaces. A comparison between our SimFWI and this separate training method is shown in Figure 2. In particular, their encoder is a pre-determined Sine kernel. Their linear layer is calculated by Ridge Regression based on the embedding from integral transforms. Then, their decoder is trained in a supervised manner with the Sine kernel encoder and frozen linear layer. However, this kernel solution has its limitations: 1) For datasets with very large variations and high-frequency components (e.g., OpenFWI Deng et al. (2022)), their performance is quite bad; 2) Since there is no explicit inverse transformation of Gaussian integral transform, the decoder still needs to be trained with supervision and cannot be shared among multiple datasets; 3) it

has very poor noise resistance. Thus, this pre-determined kernel solution may not be applicable to broader scenarios as there is no clear rule for selecting the right kernels for different situations. In contrast, we integrate self-supervised learning with the FWI problem, so that both the encoder and decoder can be pre-trained. Only a linear converter needs supervised training.

Compared to pre-determined kernel functions, properly self-supervised trained networks are more expressive, applicable to more scenarios, more noise tolerant, and can be shared across multiple datasets with different sub-surface structures. To substantiate the claims made about performance and noise resistance, we have included a comprehensive comparison in the following section. The quantitative evaluation of their performances is presented in Table 3, Section 4.2. The robustness to noise of both models is shown in Section 4.3. A comparison of the seismic and velocity latent representations obtained by our method and InvLINT is presented in the Appendix.

4 Experiments

We evaluate our approach on OpenFWI Deng et al. (2022), the first and only large-scale collection of openly accessible multi-structural seismic FWI datasets with benchmarks. We compare our method with the state-of-the-art works, including InversionNet Wu & Lin (2019), i.e., the method that jointly trains the encoder and decoder, and InvLINT Feng et al. (2022), i.e., the method that separates the encoder and decoder. We also discuss different factors that affect the performance of our method. In the Appendix, we compare the latent representation learned by our method and InvLINT, and evaluate SimFWI’s generalizability for other imaging and PDE tasks. In particular, we test it on the electromagnetic (EM) inversion task controlled by Maxwell’s equations. For interested readers, Deng et al. also provide a detailed comparison of the physics-driven method in Deng et al. (2022).

4.1 Implementation Details

Datasets. We verify our method on OpenFWI Deng et al. (2022). OpenFWI is the first open-source collection of large-scale, multi-structural benchmark datasets for data-driven seismic FWI. While real data are extremely expensive and difficult to obtain, OpenFWI is currently the largest and most comprehensive dataset. It contains 11 2D datasets with baseline, which can be divided into four groups: four datasets in the “Vel Family”, four datasets in the “Fault Family”, two datasets in the “Style Family”, and one dataset in the “Kimberlina Family”. Four datasets in the “Vel Family” are FlateVel-A/B, and CurveVel-A/B; four datasets in the “Fault Family” are FlateFault-A/B, and CurveFault-A/B; two datasets in “Style Family” are Style-A/B; and one dataset in “Kimberlina Family” is Kimberlina-CO₂. The first three families cover two versions: easy (A) and hard (B), in terms of the complexity of subsurface structures. We will use the abbreviations (e.g. FVA for FlatVel-A and CO₂ for Kimberlina-CO₂). For more details, please refer to the original paper Deng et al. (2022).

Training Details. The input seismic data are normalized to the range [-1, 1] with a log scale. We employ AdamW Loshchilov & Hutter (2018) optimizer with momentum parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$ and a weight decay of 0.05 for both self-supervision and supervision steps. In the self-supervision step, we use the same hyper-parameters and the training schedule with the original MAE paper He et al. (2022), except we change the batch size to 512 and remove the pixel normalization. We use each family together to train the MAE. Thus, in total, we trained four different models. In the supervision step, the initial learning rate is set to be 1×10^{-3} , and decayed with a cosine annealing Loshchilov & Hutter (2016). The batch size is set to 256. To make a fair comparison with the previous work, we use l_1 loss to train the linear layer. The exact network architectures are shown in Appendix. We implement our models in Pytorch and train them on 1 NVIDIA Tesla V100 GPU.

Evaluation Metrics. We apply three metrics to evaluate the geophysical properties generated by our method: MAE, MSE, and Structural Similarity (SSIM). Following the existing literature Wu & Lin (2019); Feng et al. (2022); Deng et al. (2022), MAE and MSE are employed to measure the pixel-wise error, and SSIM is to measure the perceptual similarity since velocity has highly structured information, and degradation or distortion can be easily perceived by a human. We calculate them on normalized velocity maps, i.e., MAE and MSE in the scale [-1, 1], and SSIM in the scale [0, 1].

4.2 SimFWI is Simple and Effective

Comparisons with the Joint Training Method. Table 2 shows the comparison results with InversionNet Wu & Lin (2019). The results of InversionNet are the reported benchmark in Deng et al. (2022). Compared to the jointly trained method, our SimFWI achieves comparable results on multiple datasets with only half the model size (12.3M vs. 24.4M), and only needs to supervised train the linear layer. In FlatVel-A/B, Style-B, and Kimberlina-CO₂, SimFWI even outperforms InversionNet in some metrics. The velocity maps inverted by different methods are shown in Figure 3. We can find InversionNet has a clearer boundary, while SimFWI is better at capturing the structure details in deep position (e.g., as boxed out on FlatVel-A, CurveVel-A, Style-B, and Style-A). The corresponding error map and more visualizations are provided in the Appendix for readers who might be interested. Note that, InversionNet in Style-B always outputs a strange pattern in results as boxed out in red.

Metrics	Model	FVA	FVB	CVA	CVB	FFA	FFB	CFA	CFB	SA	SB	CO2
MAE↓	SimFWI	0.0081	0.0467	0.0738	0.1820	0.0164	0.1208	0.0277	0.1791	0.0719	0.0638	0.0060
	InversionNet	0.0131	0.0351	0.0685	0.1497	0.0172	0.1055	0.0260	0.1646	0.0625	0.0689	0.0061
MSE↓	SimFWI	0.0005	0.0151	0.0188	0.1051	0.0026	0.0362	0.0061	0.0697	0.0139	0.0097	0.0017
	InversionNet	0.0004	0.0077	0.0162	0.0836	0.0018	0.0303	0.0042	0.0614	0.0105	0.0260	0.0014
SSIM↑	SimFWI	0.9888	0.9044	0.8057	0.6169	0.9701	0.6868	0.9426	0.5672	0.8423	0.7275	0.9908
	InversionNet	0.9895	0.9461	0.8074	0.6727	0.9766	0.7208	0.9566	0.6136	0.8859	0.6314	0.9872

Table 2: Quantitative results evaluated on on OpenFWI, compared with InversionNet Wu & Lin (2019), in terms of MAE, MSE, and SSIM. SimFWI achieves comparable accuracy.

Metrics	Model	FVA	FVB	CVA	CVB	FFA	FFB	CFA	CFB	SA	SB	CO2
MAE↓	SimFWI	0.0081	0.0467	0.0738	0.1820	0.0164	0.1208	0.0277	0.1791	0.0719	0.0638	0.0060
	InvLINT	0.0532	0.1621	0.0981	0.2462	0.0729	0.1522	0.0853	0.1955	0.1002	0.0835	0.0150
MSE↓	SimFWI	0.0005	0.0151	0.0188	0.1051	0.0026	0.0362	0.0061	0.0697	0.0139	0.0097	0.0017
	InvLINT	0.0085	0.0650	0.0238	0.1312	0.0190	0.0467	0.0229	0.0754	0.0209	0.0132	0.0039
SSIM↑	SimFWI	0.9888	0.9044	0.8057	0.6169	0.9701	0.6868	0.9426	0.5672	0.8423	0.7275	0.9908
	InvLINT	0.8457	0.6465	0.7355	0.4946	0.8506	0.6445	0.8204	0.5471	0.7916	0.6557	0.9760

Table 3: Quantitative results evaluated on OpenFWI, compared with InvLINT Feng et al. (2022), in terms of MAE, MSE, and SSIM. SimFWI outperforms it in terms of all three metrics. **Comparisons with the Separate Training Method.** We also compare SimFWI with InvLINT Feng et al. (2022), which also separates the encoder and decoder, and has a linear converter. Results are shown in Table 3. Compared to InvLINT, SimFWI outperforms it in terms of all three metrics. The velocity maps inverted by different methods are shown in Figure 3. The corresponding error map and more visualization results are provided in the Appendix for readers who might be interested.

We can clearly observe that InvLINT performs poorly for data with high-frequency layering locations and faults (i.e., “Vel Family” and “Fault Family”), but yields good results in smoother structures like “Style Family” and Kimberlina CO₂. This phenomenon may come from: 1) InvLINT model is very small and has limited expressive power. “Vel Family” and “Fault Family” are very diverse. It does not have enough capacity to learn all cases. 2) The Gaussian kernel cannot capture the small fault structure well, such as the interface and fault structures. 3) Their encoder uses frequency domain features. However, the high-frequency signal is mainly present in the reflected wave, which has a small amplitude. It is not easy to be captured by a frequency-domain encoder.

4.3 SimFWI has Nice Properties

In this part, we demonstrate that our SimFWI has some nice properties, including the strong generalization ability of the pre-trained encoder/decoder, good handling of noise, and solid performance on few-shot learning. There is also a correlation between our linear converter and the dataset complexity.

Generalization Ability of Encoder and Decoder. We study the generalization ability of the pre-trained encoder and decoder. In particular, we choose the seismic encoder and velocity decoder that self-supervised trained on “Fault Family”, fix it, and train the linear converter on other datasets (except the Kimberlina-CO₂, since it has different dimensions). The results are shown in Table 4.

Results show that the encoder and decoder trained on the “Fault Family” performed quite well in other datasets. They achieved impressive results, even in Style-A and Style-B, which have significantly different subsurface structures compared to the “Fault Family”. Notably, when the encoder and decoder were migrated to simpler datasets (i.e., “Vel Family”), they achieved better results than the

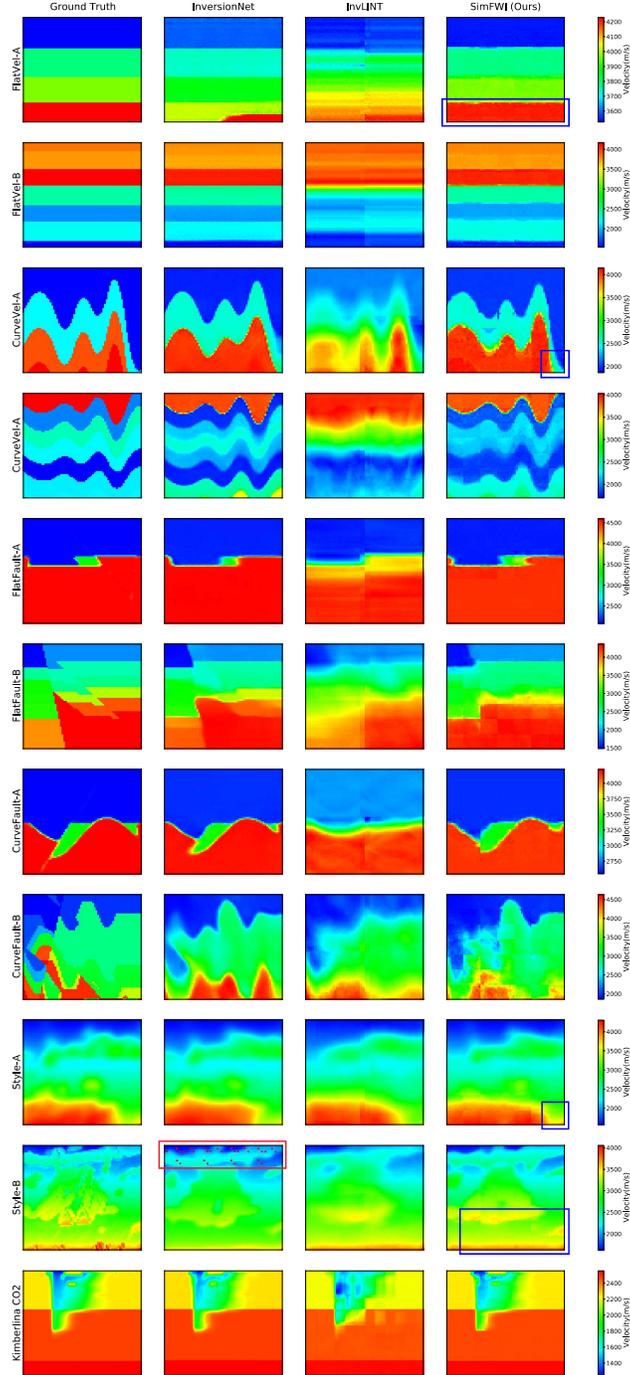


Figure 3: Results illustration of InversionNet, InvLINT and SimFWI.

results shown in Table 2 that trained exclusively on these datasets (results shown in bold). These results illustrate that 1) the latent representations obtained by self-supervision do capture the essential information in both domains and can be shared across diverse datasets, and 2) the performance of our method can be further improved by a delicate selection of self-supervision data. The current choice of using each family together in the self-supervision is to not lose generality.

To further show the generalization ability of the pre-trained encoder and decoder, and the performance improvement by picking self-supervision data, we conduct another experiment that train the encoder and decoder on cross-family datasets. In particular, CurveVel-A, FlatFault-A, and CurveFault-A are used. We test this pair of encoder and decoder in all datasets. Results are shown in Appendix.

Metrics	Model	FVA	FVB	CVA	CVB	SA	SB
MAE↓	SimFWI	0.0073	0.0570	0.0653	0.1804	0.0725	0.0646
	InversionNet	0.0131	0.0351	0.0685	0.1497	0.0625	0.0689
MSE↓	SimFWI	0.0005	0.0198	0.0159	0.1030	0.0144	0.0099
	InversionNet	0.0004	0.0077	0.0162	0.0836	0.0105	0.0260
SSIM↑	SimFWI	0.9895	0.8752	0.8192	0.6044	0.8351	0.7222
	InversionNet	0.9895	0.9461	0.8074	0.6727	0.8423	0.7275

Table 4: Generalizability of pre-trained encoder and decoder, using InversionNet as a baseline. The encoder and decoder are trained on “Fault Family”. We highlight the results that are better than the results pre-trained exclusively on these datasets.

Model	$\sigma^2=0$			$\sigma^2=1e-5$ PSNR=70.49dB			$\sigma^2=5e-5$ PSNR=63.48dB			$\sigma^2=1e-4$ PSNR=60.45dB			$\sigma^2=5e-4$ PSNR=53.39dB		
	MAE↓	MSE↓	SSIM↑	MAE↓	MSE↓	SSIM↑	MAE↓	MSE↓	SSIM↑	MAE↓	MSE↓	SSIM↑	MAE↓	MSE↓	SSIM↑
SimFWI	0.0277	0.0061	0.9426	0.0354	0.0070	0.9387	0.0508	0.0102	0.9255	0.0630	0.0139	0.9113	0.1093	0.0339	0.8308
Degradation (%)	\	\	\	-27.80	-14.75	-0.41	-83.39	-67.21	-1.81	-127.44	-127.87	-3.32	-294.58	-455.74	-11.86
InversionNet	0.0260	0.0042	0.9566	0.0332	0.0050	0.9539	0.0696	0.0133	0.9290	0.1439	0.0479	0.8830	0.4496	0.3948	0.6407
Degradation (%)	\	\	\	-27.69	-19.05	-0.28	-167.69	-216.67	-2.89	-453.46	-4.57	-7.69	-1629.23	-9300.00	-33.02
InvLINT	0.0853	0.0229	0.8204	3.1849	19.33	0.0449	7.4442	103.23	0.0172	10.16	185.87	0.0084	23.81	1033.92	0.0025
Degradation (%)	\	\	\	-3634	-84307	-94.53	-8627	-450687	-97.90	-11816	-22654	-98.98	-27807	-4514820	-99.70

Table 5: Quantitative results of different models on CurveFault-A with noisy seismic input. Gaussian noise with different variance σ^2 is added to seismic data during testing.

Good handling of noise. We provide the quantitative results of the robustness test. In particular, we add Gaussian Noise with different variances to the input seismic data during testing. Table 5 shows the performance on CurveFault-A for clean and noisy data. We also include the noise’s variance (σ^2) and average peak-to-noise ratio (PSNR) in the table. The PSNR of a sample is defined as

$$\text{PSNR} = 10 \log_{10} \frac{(p_{max} - p_{min})^2}{\ell_2(p - p')}, \quad (2)$$

where p_{max} and p_{min} denote the maximum and minimum possible values of the seismic data in a dataset, p is the clean seismic data, and p' is the noisy data.

Compared to other models, our SimFWI is the most robust one to noise. The robustness of SimFWI shows in two aspects. First, its performance degradation on noisy data is smaller than others. Second, when the noise’s variance is large ($\sigma^2 \geq 5e-5$), our method outperforms InversionNet. As expected, InvLINT is extremely sensitive to the noise, since it only uses a Fourier transform as its encoder.

Strong Performance on Few-Shot Learning. One of the most important benefits of our method is it does not need paired data to train its encoder and decoder. Thus, we test SimFWI on the Few-shot learning situation, where only a limited number of paired data exists, and compare it with the jointly trained method, InversionNet. We choose three datasets as examples and test the situation that only 1/5 or 1/10 paired data can be used in supervised learning. MAE results are reported in Table 6. We can see that on all three datasets, our approach’s performance is a little worse than InversionNet when using the full data. However, when the paired data becomes less, our model outperforms InversionNet. This implies that our SimFWI has a strong performance on few-shot learning.

Correlation between the linear layer and datasets’ complexity. By simplifying the image-to-image translation problem to a linear problem, our model is easy to analyze. With only a linear converter trained in a supervised manner, we can conduct singular value decomposition analysis.

Dataset	Model	Ratio=1	Ratio=1/5	Ratio=1/10
CurveVel-A	SimFWI	0.0738	0.1108	0.1284
	InversionNet	0.0685	0.1113	0.1335
CurveFault-A	SimFWI	0.0277	0.0500	0.0683
	InversionNet	0.0260	0.0589	0.0843
Style-A	SimFWI	0.0719	0.0917	0.1030
	InversionNet	0.0625	0.0942	0.1046

Table 6: MAE results using partial datasets. Ratio indicates the proportion of datasets used.

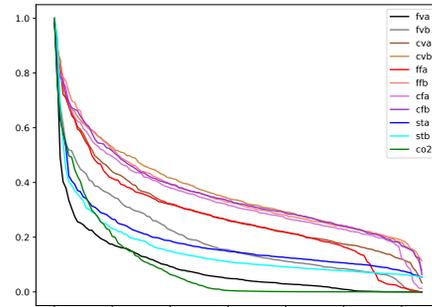


Table 7: Normalized Singular Values.

The results are shown in Figure 7. Since the singular value varies greatly in different datasets, we normalize it by dividing it by its maximum value and trunk it at 128 dim, which is the bottleneck dimension.

We can clearly observe that there is a strong correlation between the singular values and the dataset complexity. Generally speaking, CurveFault-B, FlatFault-B, Style-A, Style-B, and CurveFault-B have the most complex velocity map among all datasets. As we can see, their singular values are much slower to fall. On the other hand, FlatVel-A and Kimberlina-CO₂ are the simplest datasets, which are also reflected in their singular values. We also provide the original singular Value in Appendix Figure 4 for reference. All these results prove that our linear converter correlates to the datasets’ complexity. For readers who might be interested, Deng et al. (2022) provides a more detailed analysis of datasets’ complexity among all 11 datasets.

4.4 Ablation Test

In this part, we will test how the rank of the linear converter will influence the performance and demonstrates the influence of the local linear relationship between two latent spaces. In the Appendix, we also show the comparison between using the masked autoencoder and the classical autoencoder, and test the performance with several different non-linear converters.

Rank of Linear. We evaluate performances over five different numbers of ranks of the linear converter, varying from 32 to 128. The quantitative results are shown in Table 8. Results indicate that increasing the rank makes the model much larger, but the growth of the results is limited. On the other hand, decreasing the model’s rank also does not reduce its capacity a lot but results in a smaller number of parameters. This allows for the balance of performance and computational cost based on specific requirements and available resources, highlighting the flexibility of our model.

Dataset	Dim	#Param	MAE↓	MSE↓	SSIM↑
CurveFault-A	128*	12.3M	0.0277	0.0061	0.9426
	512	26.0M	0.0271	0.0058	0.9441
	256	16.8M	0.0274	0.0059	0.9434
	64	10.0M	0.0280	0.0064	0.9392
	32	8.9M	0.0304	0.0075	0.9300

Table 8: Quantitative results of different bottleneck dimensions, and the corresponding number of parameters. As a reference, InversionNet has 24.4M parameters. (*) indicates the default option.

Local Linear Relationship. We demonstrate the property we found that each dataset shows linear relation locally, and there is a piece-wise linear relation globally over multiple datasets. In particular, we let datasets in each family share not only the encoder and decoder but also the linear converter. In other words, we use all datasets in each family to train the linear converter. We report the results and performance change in Table 9. In the table, we highlight the improvement of the results after sharing the linear converter. It is quite interesting that, generally, the datasets with a more complex subsurface structure show a performance improvement. In contrast, simpler datasets’ performance drops a lot. The results come from the fact that a complex dataset covers a larger range in the latent space. The scope of simple datasets is covered by those complex ones in the same family. Thus, with more data to use, SimFWI achieves better results on complex datasets. But, for simple datasets, out-of-distribution data make the learning results deviate substantially from their local linear relationship.

Metrics	Model	FlatVel-A	FlatVel-B	CurveVel-A	CurveVel-B	FlatFault-A	FlatFault-B	CurveFault-A	CurveFault-B	Style-A	Style-B
MAE↓	Original	0.0081	0.0467	0.0738	0.1820	0.0164	0.1208	0.0277	0.1791	0.0719	0.0638
	Sharing Linear	0.0191	0.0545	0.0761	0.1709	0.0306	0.1198	0.0421	0.1697	0.0699	0.0636
MSE↓	Original	0.0005	0.0151	0.0188	0.1051	0.0026	0.0362	0.0061	0.0697	0.0139	0.0097
	Sharing Linear	0.0015	0.0161	0.0193	0.0963	0.0054	0.0339	0.0093	0.063	0.013	0.0094
SSIM↑	Original	0.9888	0.9044	0.8057	0.6169	0.9701	0.6868	0.9426	0.5672	0.8423	0.7275
	Sharing Linear	0.9633	0.8827	0.8007	0.6326	0.9411	0.6913	0.9096	0.5873	0.8512	0.7248

Table 9: Quantitative results of sharing linear converter over multiple datasets, compared with original results. Both the encoder/decoder and the linear layer are shared across each dataset family. We highlight the improvement of the results after sharing the linear converter.

5 Discussion

While our proposed method has shown promising results, there are still some limitations that need to be further addressed. Firstly, the linear relation in the latent space is currently only piece-wise linear.

Although we have demonstrated the effectiveness of sharing the encoder and decoder across datasets, the linear converter is still limited to a single dataset. Thus, there is a need to explore more advanced methods for learning better representations that make the latent-space representations closer to the physical nature and enable the linear relationship to be applicable in broader scenarios.

Secondly, it is hard to train our model from scratch exclusively on real data due to the lack of labeled real data in subsurface geophysics. This is not just a limitation of our work, but a common limitation in the seismic inversion community. The concept of “Sim2Real” is a well-received technique to transfer knowledge learned in simulation to real data James et al. (2019). To mitigate the gap between simulation and real scenarios, we have tested our model in velocity maps that yield physically realistic subsurface structures, i.e. Style-A and Style-B Feng et al. (2021). Additionally, we have imposed noise to simulate more realistic measurement procedures. Our method demonstrated promising performance in both scenarios of realistic subsurface data and noisy seismic measurements. We will explore how to train the converter with purely unpaired data and further mitigate the knowledge gap when applying our proposed model to real data in our future work.

6 Conclusion

In this paper, we present a new framework, SimFWI, that simplifies the mapping between seismic data and velocity maps into a linear problem via domain-independent self-supervised learning. We decouple two domains of FWI and train the encoder and decoder separately in their own domain, with two masked autoencoders. We observed a linear correlation between the two latent spaces, meaning that the self-supervised encoder and decoder can be frozen, and a linear converter can be learned to connect them from the paired seismic data and velocity map. This framework allows a better understanding of the relationship among multiple FWI datasets with different subsurface structures. In experiments, SimFWI achieved comparable performance, with half the model size, and showed solid performance in a few-shot situation and robustness test.

References

- Alumbaugh, D., Commer, M., Crandall, D., Gasperikova, E., Feng, S., Harbert, W., Li, Y., Lin, Y., Manthila Samarasinghe, S., and Yang, X. Development of a multi-scale synthetic data set for the testing of subsurface CO₂ storage monitoring strategies. In American Geophysical Union (AGU), 2021.
- Araya-Polo, M., Jennings, J., Adler, A., and Dahlke, T. Deep-learning tomography. The Leading Edge, 37(1):58–66, 2018.
- Chen, D., Tachella, J., and Davies, M. E. Equivariant imaging: Learning beyond the range space. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 4379–4388, 2021.
- Deng, C., Feng, S., Wang, H., Zhang, X., Jin, P., Feng, Y., Zeng, Q., Chen, Y., and Lin, Y. Openfwi: Large-scale multi-structural benchmark datasets for full waveform inversion. 2022.
- Feng, S., Fu, L., Feng, Z., and Schuster, G. T. Multiscale phase inversion for vertical transverse isotropic media. Geophysical Prospecting, 69(8-9):1634–1649, 2021.
- Feng, Y., Chen, Y., Feng, S., Jin, P., Liu, Z., and Lin, Y. An intriguing property of geophysics inversion. The Thirty-ninth International Conference on Machine Learning, 2022.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16000–16009, 2022.
- James, S., Wohlhart, P., Kalakrishnan, M., Kalashnikov, D., Irpan, A., Ibarz, J., Levine, S., Hadsell, R., and Bousmalis, K. Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12627–12637, 2019.

- Jin, P., Zhang, X., Chen, Y., Huang, S. X., Liu, Z., and Lin, Y. Unsupervised learning of full-waveform inversion: Connecting CNN and partial differential equation in a loop. In Proceedings of the Tenth International Conference on Learning Representations (ICLR), 2022.
- Lin, Y., Theiler, J., and Wohlberg, B. Physics-guided data-driven seismic inversion: Recent progress and future opportunities in full waveform inversion. IEEE Signal Processing Magazine, 40:115–133, 2023.
- Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983, 2016.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In Sixth International Conference on Learning Representations (ICLR), 2018.
- Sun, J., Innanen, K. A., and Huang, C. Physics-guided deep learning for seismic inversion with hybrid training and uncertainty analysis. Geophysics, 86(3):R303–R317, 2021.
- Wu, Y. and Lin, Y. InversionNet: An efficient and accurate data-driven full waveform inversion. IEEE Transactions on Computational Imaging, 6:419–433, 2019.
- Zeng, Q., Feng, S., Wohlberg, B., and Lin, Y. Inversionnet3d: Efficient and scalable learning for 3-d full-waveform inversion. IEEE Transactions on Geoscience and Remote Sensing, 60:1–16, 2021.
- Zhang, Z., Wu, Y., Zhou, Z., and Lin, Y. Velocitygan: Subsurface velocity image estimation using conditional adversarial networks. In 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), pp. 705–714. IEEE, 2019.

A Appendix

A.1 Architecture

The exact transformer architectures and layer dimensions of the seismic and velocity autoencoders are provided in Table 10. We will add this table to the revised paper. For all the datasets except Kimberlina, the size of seismic data is 1000×70 and the size of velocity maps is 70×70 . We choose the patch size 100×10 for seismic data and 10×10 for velocity maps. Thus, the latent dimension of seismic data is 132×70 , and the latent dimension of velocity maps is 516×49 . For Kimberlina, the patch size of seismic data is 250×10 , and of velocity maps is 20×40 . The latent dimension of seismic data is 132×50 , and the latent dimension of velocity maps is 516×70 . The rank of the linear converter is set to 128. The mask ratio for training MAE is set to 0.75.

Model	#Layers	Embedded Dim	MLP Dim	#Heads
Seismic Encoder	2	132	528	12
Seismic Decoder	2	512	144	16
Velocity Encoder	3	516	2064	12
Velocity Decoder	2	512	2064	16

Table 10: Details of seismic and velocity autoencoders

A.2 Generalizability

Dataset	MAE↓	MSE↓	SSIM↑
CurveVel-A*	0.0634	0.0155	0.8267
FlatFault-A*	0.0166	0.0026	0.9698
CurveFault-A*	0.0271	0.006	0.9434
FlatVel-A	0.0072	0.0004	0.9912
FlatVel-B	0.0552	0.0179	0.8783
CurveVel-B	0.1754	0.0981	0.6157
FlatFault-B	0.1260	0.0381	0.6734
CurveFault-B	0.1837	0.0711	0.5590
Style-A	0.0744	0.0146	0.8311
Style-B	0.0653	0.0102	0.7175

Table 11: Quantitative results of the generalization ability of pre-trained encoder and decoder. The encoder and decoder are trained across datasets’ families. (*) indicates the datasets used to train the encoder and decoder.

A.3 Ablation Test

In this part, we show the comparison between using masked autoencoder and autoencoder as self-supervised learners; and test the performance of several different non-linear converters.

MAE v.s. Autoencoder. We conducted another experiment that use autoencoders (i.e., mask ratio equals zero.) with the same architecture as the self-supervised training model. We pre-trained the model on “Fault Family”; and trained the linear converter and validated it on CurveFault-A as an example. The reconstruction and inversion results are shown in Table 12. As demonstrated, a simple autoencoder cannot capture the important information that is necessary for both reconstructing and connecting to another domain. If we simply consider both seismic data and velocity maps as pure images and ignore the physical meaning behind them, the autoencoder would learn too many shortcuts which are only useful to reconstruct the image but lost the essential information reflecting its physical properties. This is because seismic data and velocity maps not as diverse as natural images. On the other hand, if a model can embed the essential underlying physics information of these two quantities, it will naturally enhance the generalization ability.

Non-Linear Converter. We evaluate networks with a more complicated nonlinear converter on CurveFault-A. We tested four different settings: 1) a two-layer MLP; 2) a two-piece Maxout layer; 3) a two-layer U-Net; and 4) a four-layer U-Net. The results are provided in Table 13. From the results, we can see that 1) a simple nonlinear mapping (e.g., two-layer MLP or U-Net) has no positive

Model	MAE↓	MSE↓	SSIM↑	Seismic	Velocity
				Pre-training MAE↓	Pre-training MAE↓
Masked Autoencoder	0.0277	0.0061	0.9426	0.1703	0.0410
Autoencoder	0.0614	0.0174	0.8302	0.0008	0.0005

Table 12: Comparison between different pre-training strategies on CurveFault-A. In addition to the quantitative results of inversion, the mean absolute reconstruction errors (with masks) of the pre-trained models (Columns 5 & 6) are also reported.

effect on final performance; and 2) a piece-wise linear mapping (Maxout) or a much more complex nonlinear mapping (four-layer U-Net) can only provide limited improvement. These results are consistent with our conclusion of a near-linear relationship.

Model	MAE↓	MSE↓	SSIM↑
Linear	0.0277	0.0061	0.9426
Two-Layer MLP	0.0280	0.0064	0.9433
Two-Pieces Maxout	0.0260	0.0057	0.9472
2-Layer U-Net	0.0285	0.0062	0.9414
4-Layer U-Net	0.0259	0.0056	0.9465

Table 13: Quantitative results on CurveFault-A with different nonlinear converters.

A.4 Singular Value Decomposition

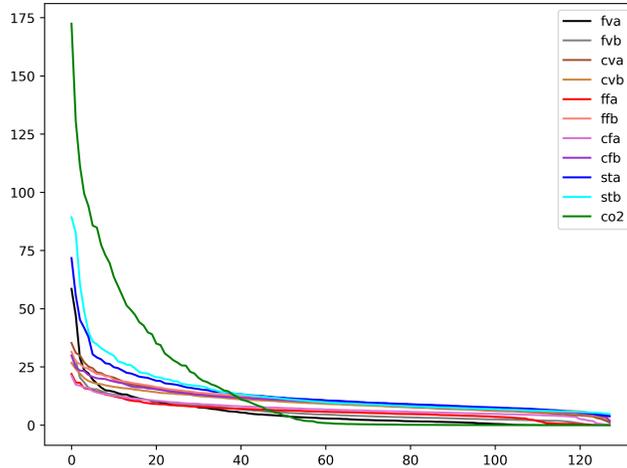


Figure 4: Original Results of Singular Value Decomposition on different datasets.

A.5 Comparing the latent representations of SimFWI and InvLINT.

To further analyze the relation between our SimFWI and InvLINT, in this part, we compare the latent representations of seismic data and velocity maps obtained by our method to those obtained by InvLINT. First, We conducted experiments on CurveFault-A that use a sine kernel from InvLINT as the encoder and use our pre-trained decoder to construct the inversion network, respectively. The converter is still linear. The results are shown in Table 14. These results show that using the latent seismic representation from the sine kernel is difficult to regress the latent velocity representation from our method.

To further compare the latent representations, we use one latent representation to predict another with linear regression, for seismic data and velocity maps respectively. We report the coefficient of determination (R^2 score) in Table 15.

These show that our latent space with a higher dimension contains more information. As a preliminary comparison, we can roughly conclude that their latent space is a linear subspace of our latent space.

Model	MAE↓	MSE↓	SSIM↑
SimFWI	0.0277	0.0061	0.9426
Sine Kernel Encoder	0.0426	0.0093	0.9233

Table 14: Comparison between latent representations of seismic data obtained by SimFWI and InvLINT on CurveFault-A.

Variable	Source	Target	R^2
Seismic	SimFWI	InvLINT	0.9869
	InvLINT	SimFWI	0.6700
Velocity	SimFWI	InvLINT	0.9996
	InvLINT	SimFWI	0.4871

Table 15: Predicting the target latent representations from the source latent representations with linear regression.

A.6 Generalizability of the SimFWI for the other inversion task.

We conducted an experiment using different PDEs on the Kimberlina-Reservoir dataset Alumbaugh et al. (2021); Feng et al. (2022). In this dataset, the task is to recover the subsurface conductivity from electromagnetic (EM) measurements acquired on the surface. The governing equations here are Maxwell’s Equations

$$\begin{aligned}\sigma \mathbf{E} - \nabla \times \mathbf{H} &= -\mathbf{J}, \\ \nabla \times \mathbf{E} + i\omega\mu_0 \mathbf{H} &= -\mathbf{M},\end{aligned}$$

where \mathbf{E} and \mathbf{H} are the electric and magnetic fields. \mathbf{J} and \mathbf{M} are the electric and magnetic sources. σ is the electrical conductivity and $\mu_0 = 4\pi \times 10^{-7} \Omega \cdot s/m$ is the magnetic permeability of free space. We compared the results of our SimFWI model with those reported in InvLINT (Feng et al., 2022), and presented the results in Table 16. Note that, to maintain consistency with InvLINT, the MAE and MSE reported below were calculated after denormalizing to the original range of $[0, 0.65]$. For all other results presented in our paper, the MAE and MSE were calculated in the normalized range of $[-1, 1]$. We observe that our proposed SimFWI yields better performance than those obtained using InvLINT and InversionNet.

Dataset	Model	MAE↓	MSE↓	SSIM↑
Kimberlina-Reservoir	SimFWI	0.00438	0.000192	0.9700
	InversionNet	0.01330	0.000855	0.9175
	InvLINT	0.00703	0.000537	0.9370

Table 16: Quantitative results on Kimberlina-Reservoir. MAE and MSE are calculated after denormalizing to their original range (i.e., $[0, 0.65]$).

A.7 Visualizations

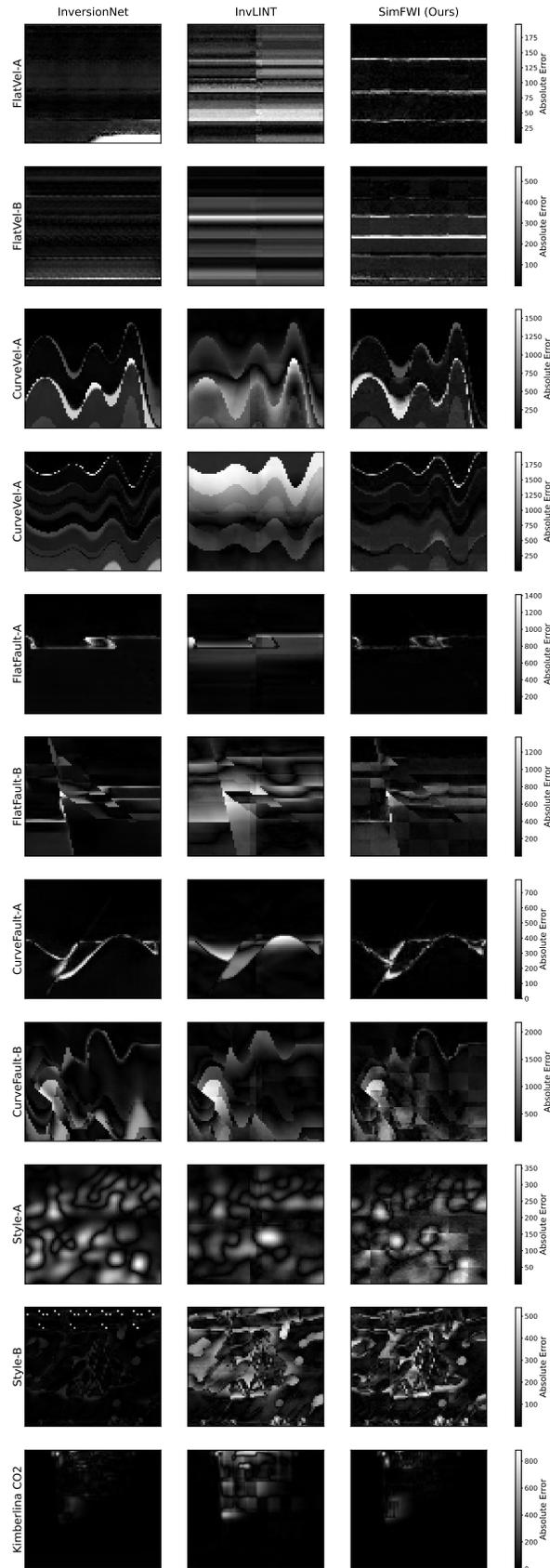


Figure 5: Illustration of absolute error map on OpenFWI, compared, SimFWI, InversionNet Wu & Lin (2019) and InvLINT Feng et al. (2022) to the ground truth.

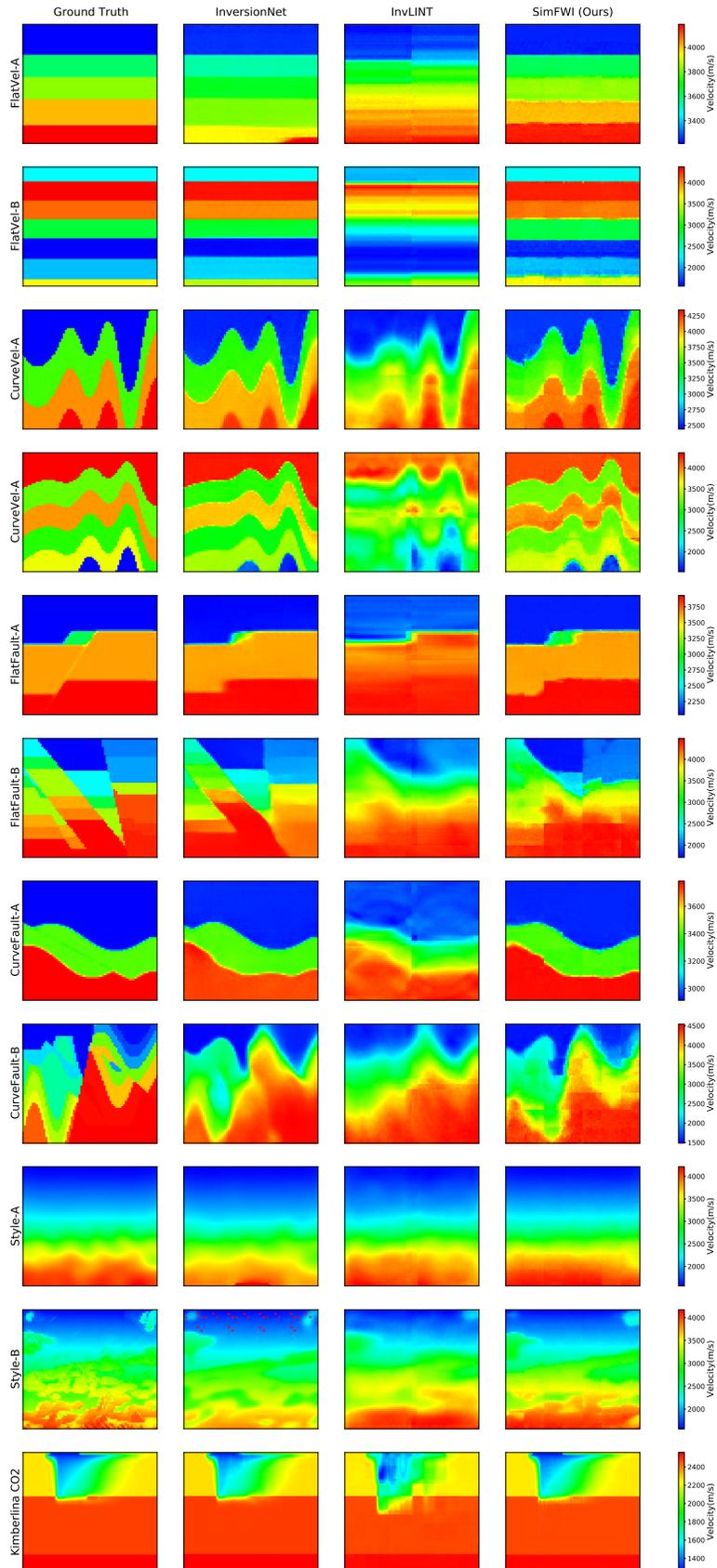


Figure 6: Illustration of results evaluated on OpenFWI, compared with InversionNet Wu & Lin (2019) and InvLINT Feng et al. (2022).

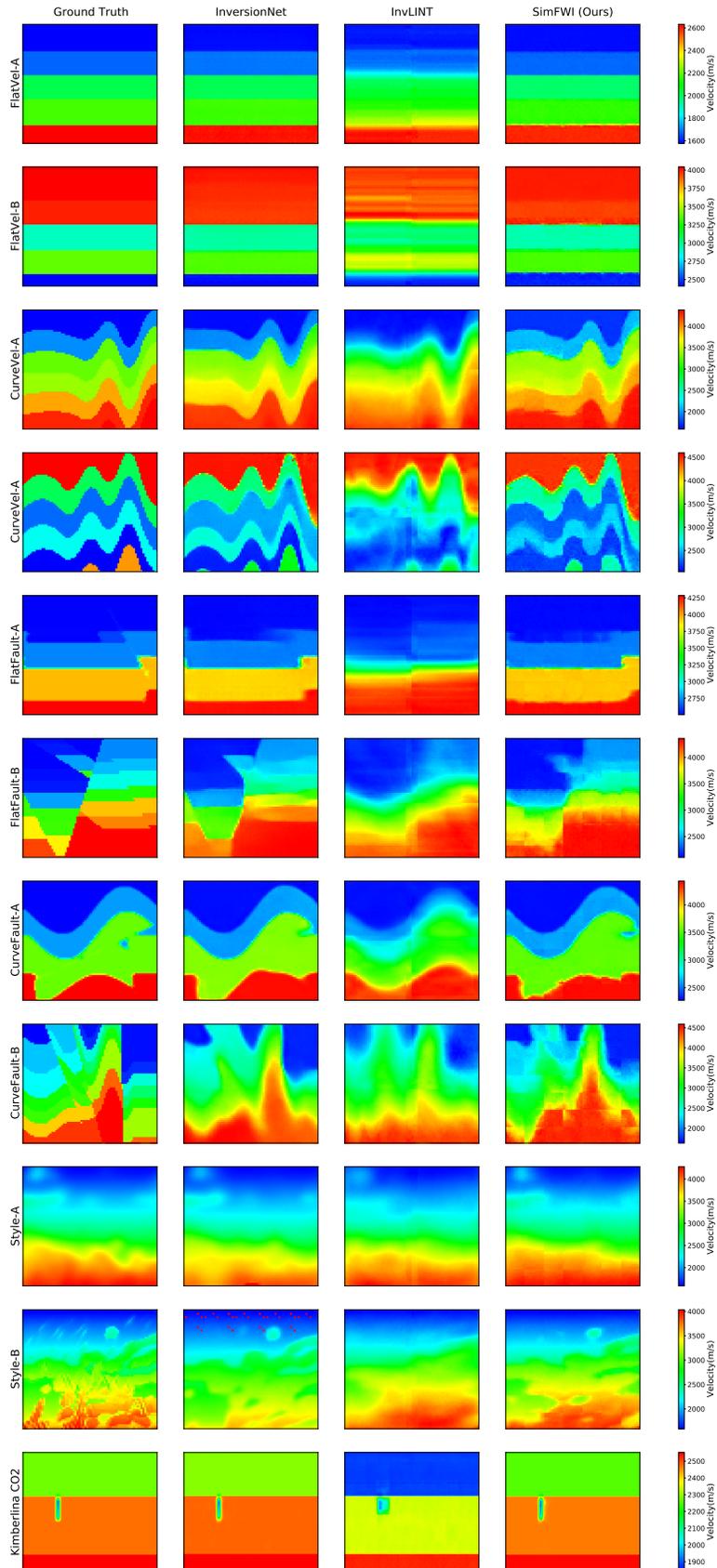


Figure 7: Illustration of results evaluated on OpenFWI, compared with InversionNet Wu & Lin (2019) and InvLINT Feng et al. (2022).