# A deep learning perspective of the forward and inverse problems in exploration geophysics

Jian Sun*, Zhan Niu, Kris Innanen, Junxiao Li*, Daniel Trad

Nov 30, 2018

Banff, AB

# Outline

➢ Motivation

# Outline

➢ Motivation

➢ The forward problem: a deep learning perspective

■ Forward modeling of wave propagation

■ Recurrent Neural Network (RNN)

➢ Motivation

➢ The forward problem: a deep learning perspective

- ■ Forward modeling of wave propagation
- ■ Recurrent Neural Network (RNN)

➢ The inverse problem: a deep learning perspective

- ■ The gradient derivation in a RNN framework
- ■ Connections with FWI?

# Outline

➢ Motivation
➢ The forward problem: a deep learning perspective
  ▪ Forward modeling of wave propagation
  ▪ Recurrent Neural Network (RNN)
➢ The inverse problem: a deep learning perspective
  ▪ The gradient derivation in a RNN framework
  ▪ Connections with FWI?
➢ Numerical analysis & tuning of hyperparameters
  ▪ Best learning rate selection for gradient-based algorithms
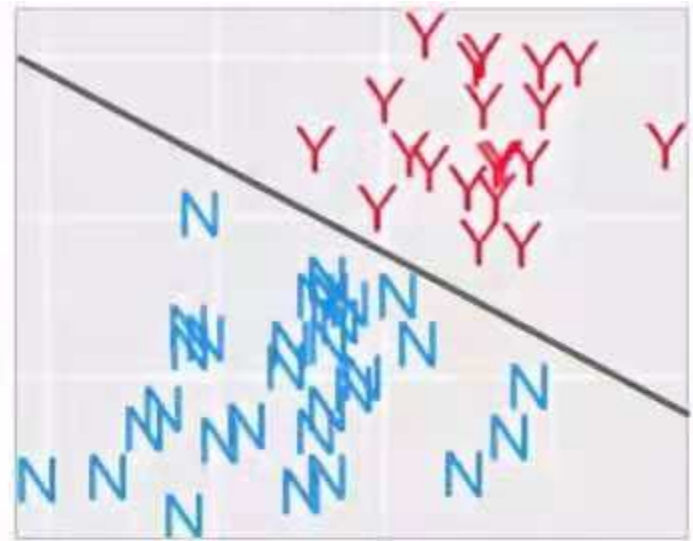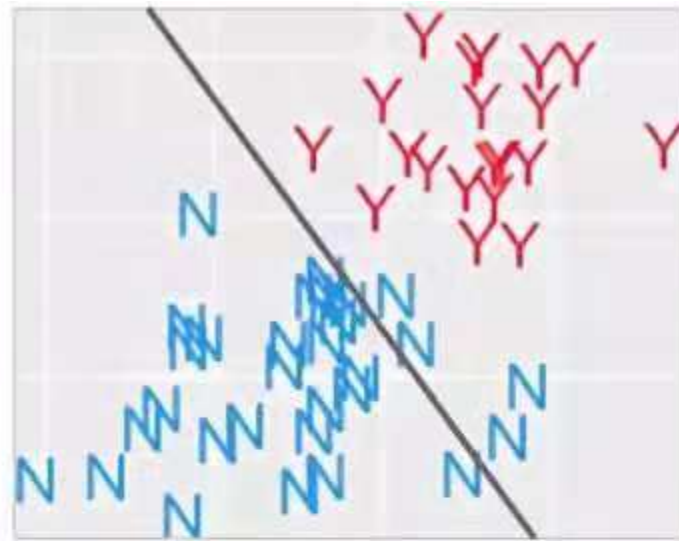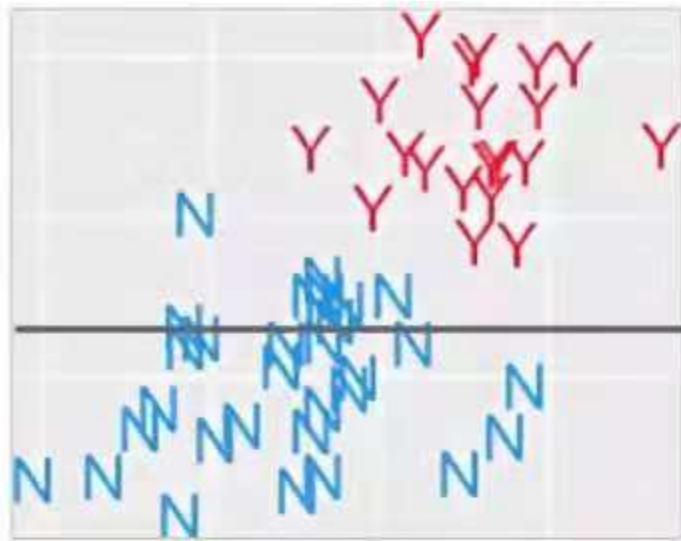  ▪ Comparisons (GD, Momentum, Adagrad, RMSprop, Adam, CG, L-BFGS)

# Outline

➢ **Motivation**

➢ **The forward problem: a deep learning perspective**
  - Forward modeling of wave propagation
  - Recurrent Neural Network (RNN)

➢ **The inverse problem: a deep learning perspective**
  - The gradient derivation in a RNN framework
  - Connections with FWI?

➢ **Numerical analysis & tuning of hyperparameters**
  - Best learning rate selection for gradient-based algorithms
  - Comparisons (GD, Momentum, Adagrad, RMSprop, Adam, CG, L-BFGS)

➢ **Synthetic test on Marmousi**

# Outline

- Motivation
- The forward problem: a deep learning perspective
    - Forward modeling of wave propagation
    - Recurrent Neural Network (RNN)
- The inverse problem: a deep learning perspective
    - The gradient derivation in a RNN framework
    - Connections with FWI?
- Numerical analysis & tuning of hyperparameters
    - Best learning rate selection for gradient-based algorithms
    - Comparisons (GD, Momentum, Adagrad, RMSprop, Adam, CG, L-BFGS)
- Synthetic test on Marmousi
- Conclusions & Future works

# Motivation

➢ Deep learning (DL):
  ✓ Widely used: speech recognition, computer vision, auto-driving, machine translation, medical imaging, etc.
  ✓ Fast: only trained once.
  ✗ Large and well-labelled training datasets.
  ✗ Not theory-guided (data-determined).

➢ Full waveform inversion (FWI):
  ✓ Theory-guided: small datasets.
  ✓ Full wavefield information used.
  ✗ Computationally expensive.
  ✗ Cycle-skipping.

- Forward modeling of wave propagation

$$\nabla^2 \mathbf{u}(\mathbf{r}, t) = \frac{1}{v^2(\mathbf{r})} \frac{\partial^2 \mathbf{u}(\mathbf{r}, t)}{\partial t^2} + s(\mathbf{r}, t)\delta(\mathbf{r} - \mathbf{r_s})$$

- The second-order finite-difference:

$$\mathbf{u}(\mathbf{r}, t + \Delta t) = v^2(\mathbf{r})\Delta t^2 \left[\nabla^2 \mathbf{u}(\mathbf{r}, t) - s(\mathbf{r}, t)\delta(\mathbf{r} - \mathbf{r_s})\right] + 2\mathbf{u}(\mathbf{r}, t) - \mathbf{u}(\mathbf{r}, t - \Delta t)$$

- Forward modeling of wave propagation

$$\nabla^2 \mathbf{u}(\mathbf{r}, t) = \frac{1}{v^2(\mathbf{r})} \frac{\partial^2 \mathbf{u}(\mathbf{r}, t)}{\partial t^2} + s(\mathbf{r}, t)\delta(\mathbf{r} - \mathbf{r_s})$$

- The second-order finite-difference:

$$\mathbf{u}(\mathbf{r}, t + \Delta t) = v^2(\mathbf{r})\Delta t^2 \left[ \nabla^2 \mathbf{u}(\mathbf{r}, t) - s(\mathbf{r}, t)\delta(\mathbf{r} - \mathbf{r_s}) \right] + 2\mathbf{u}(\mathbf{r}, t) - \mathbf{u}(\mathbf{r}, t - \Delta t)$$

$u_t$

$u_{t-\Delta t}$

FD Operator

$u_{t+\Delta t}$

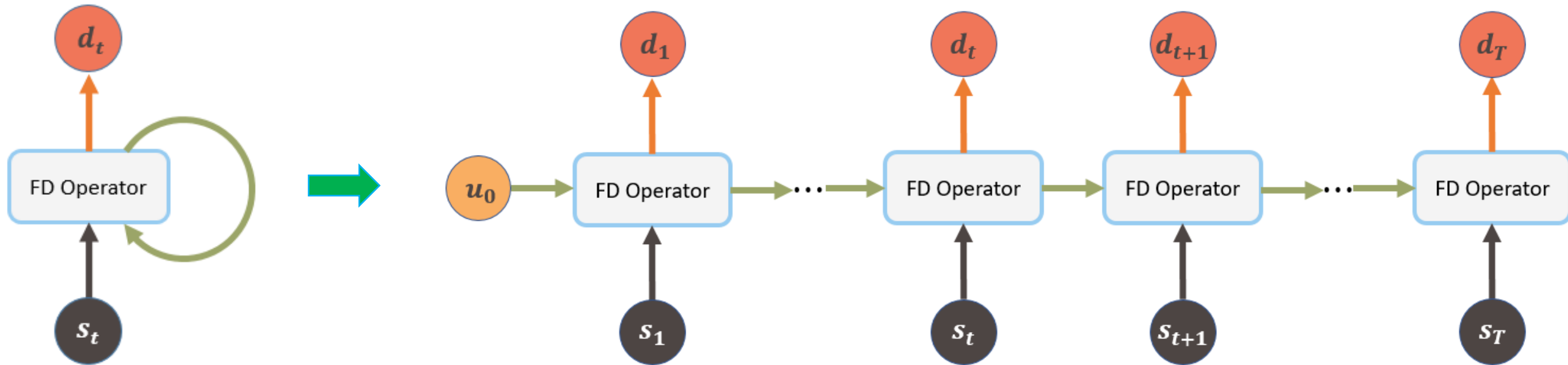# Recurrent Neural Network (RNN)

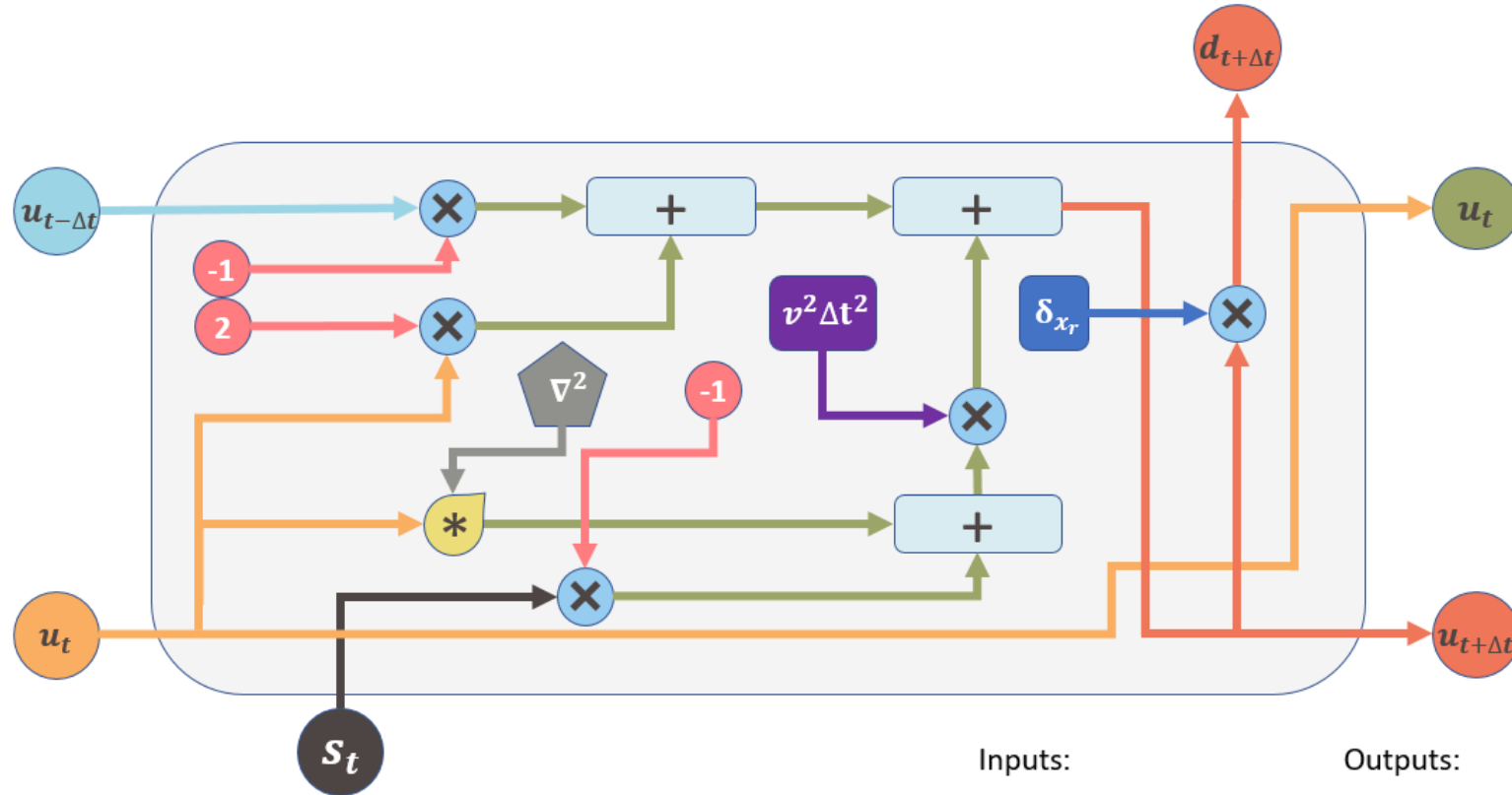

one to one | one to many | many to one | many to many | many to many

# Recurrent Neural Network (RNN)

# Recurrent Neural Network (RNN)

- The gradient derivation in a RNN framework



> Objective function: $J(v) = \dfrac{1}{2n_s} \sum\limits_{r_s} \sum\limits_{r_g} \sum\limits_{t} (\mathbf{d_t} - \tilde{\mathbf{d}}_\mathbf{t})^2 = \dfrac{1}{2n_s} \sum\limits_{r_s} \sum\limits_{r_g} \sum\limits_{t} (\mathbf{d_t} - \boldsymbol{\delta}_{\mathbf{r_g}} \tilde{\mathbf{u}}_\mathbf{t})^2$

> Gradient: $\dfrac{\partial J}{\partial v} = \sum\limits_{\mathbf{t}}^{\mathbf{T}} \left[ \dfrac{\partial J}{\partial \tilde{\boldsymbol{u}}_t} \right] \dfrac{\partial \tilde{\boldsymbol{u}}_t}{\partial v}$
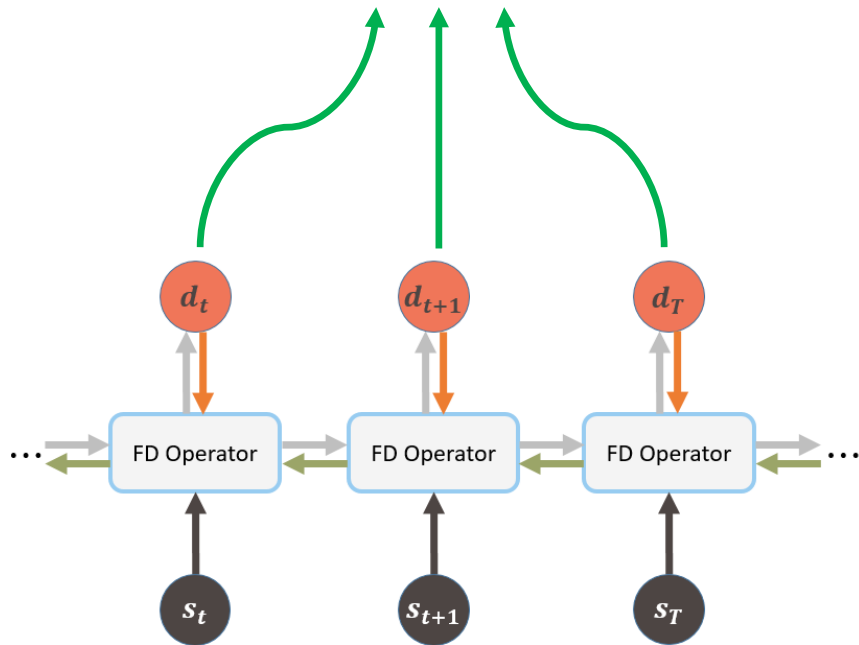
- The gradient derivation in a RNN framework

$$J(\boldsymbol{v}) = \frac{1}{2n_s} \sum_{\mathbf{r_s}} \sum_{\mathbf{r_g}} \sum_{\mathbf{t}} (\mathbf{d_t} - \boldsymbol{\delta_{r_g}} \tilde{\mathbf{u}}_\mathbf{t})^2$$



➢ Gradient at time-step t:
$$\left[\frac{\partial J}{\partial \boldsymbol{v}}\right]_t = \left[\frac{\partial J}{\partial \tilde{\boldsymbol{u}}_t}\right] \frac{\partial \tilde{\boldsymbol{u}}_t}{\partial \boldsymbol{v}}$$

$$\left[\frac{\partial J}{\partial \tilde{u}_t}\right] = \left[\frac{\partial J}{\partial \tilde{u}_{t+2}}\right] \frac{\partial \tilde{u}_{t+2}}{\partial \tilde{u}_t} + \left[\frac{\partial J}{\partial \tilde{u}_{t+1}}\right] \frac{\partial \tilde{u}_{t+1}}{\partial \tilde{u}_t} + \frac{\partial J}{\partial \tilde{u}_t}$$

$$= \boldsymbol{v}^2 \Delta t^2 \left( \nabla^2 \left[\frac{\partial J}{\partial \tilde{u}_{t+1}}\right] - \frac{1}{n_s \boldsymbol{v}^2 \Delta t^2} \sum_{\mathbf{r_s}} \sum_{\mathbf{r_g}} \boldsymbol{\delta} \mathbf{d_t} \right)$$

$$+ 2 \left[\frac{\partial J}{\partial \tilde{u}_{t+1}}\right] - \left[\frac{\partial J}{\partial \tilde{u}_{t+2}}\right]$$

$$\frac{\partial \tilde{u}_t}{\partial \boldsymbol{v}} = \frac{2\Delta t^2}{\boldsymbol{v}} \boldsymbol{v}^2 (\nabla^2 \tilde{\mathbf{u}}_{\mathbf{t-1}} - \boldsymbol{s}_{\mathbf{t-1}}) \approx \frac{2\Delta t^2}{\boldsymbol{v}} \frac{\partial^2 \tilde{u}_t}{\partial t^2}$$

- ■ Connections with FWI?

$$J(\boldsymbol{v}) = \frac{1}{2n_s} \sum_{\mathbf{r_s}} \sum_{\mathbf{r_g}} \sum_{\mathbf{t}} (\mathbf{d_t} - \boldsymbol{\delta_{r_g}} \tilde{\mathbf{u}}_\mathbf{t})^2$$



➢ Gradient at time-step t:

$$\boldsymbol{g}_t = \left[\frac{\partial J}{\partial \boldsymbol{v}}\right]_t$$

$$= BP\left(-\frac{1}{n_s \boldsymbol{v}^2 \Delta t^2} \sum_{\mathbf{r_s}} \sum_{\mathbf{r_g}} \boldsymbol{\delta}\mathbf{d_t}\right) \frac{2\Delta t^2}{\boldsymbol{v}} \frac{\partial^2 \tilde{\boldsymbol{u}}_{t-1}}{\partial t^2} \qquad \text{⬅ RNN}$$

$$\approx BP\left(-\frac{1}{n_s} \sum_{\mathbf{r_s}} \sum_{\mathbf{r_g}} \boldsymbol{\delta}\mathbf{d_t}\right) \frac{2}{\boldsymbol{v}^3} \frac{\partial^2 \tilde{\boldsymbol{u}}_t}{\partial t^2} \qquad \text{⬅ FWI}$$

✓ Which means, training this self-designed RNN is approximately equivalent to the FWI process. In other words, FWI is also a special case of machine learning task.

# Numerical analysis & tuning of hyperparameters

- Find the best learning rate ranges for gradient-based algorithms including GD, Momentum, Adagrad, RMSprop, Adam
- Several inter-comparisons

$$\boldsymbol{v}_k = \boldsymbol{v}_{k-1} - \alpha \cdot \boldsymbol{g}_k$$

$\alpha \in (0,1]$

$$\boldsymbol{m}_k = \beta \cdot \boldsymbol{m}_{k-1} + (1 - \beta) \cdot \boldsymbol{g}_k$$

$$\tilde{\boldsymbol{m}}_k = \boldsymbol{m}_k / (1 - \beta^k)$$

$$\boldsymbol{v}_k = \boldsymbol{v}_{k-1} - \alpha \cdot \tilde{\boldsymbol{m}}_k$$

$$\alpha \in (0,1]$$

$$\boldsymbol{m}_k = \beta \cdot \boldsymbol{m}_{k-1} + (1 - \beta) \cdot \boldsymbol{g}_k$$

$$\tilde{\boldsymbol{m}}_k = \boldsymbol{m}_k / (1 - \beta^k)$$

$$\boldsymbol{v}_k = \boldsymbol{v}_{k-1} - \alpha \cdot \tilde{\boldsymbol{m}}_k$$

$$\alpha \in (0,1]$$

$$G_{k,ii} = G_{k-1,ii} + g_{k,i}^2$$

$$v_{k,i} = v_{k-1,i} - \boxed{\frac{\alpha}{\sqrt{G_{k,ii} + \epsilon}}} \cdot g_{k,i}$$
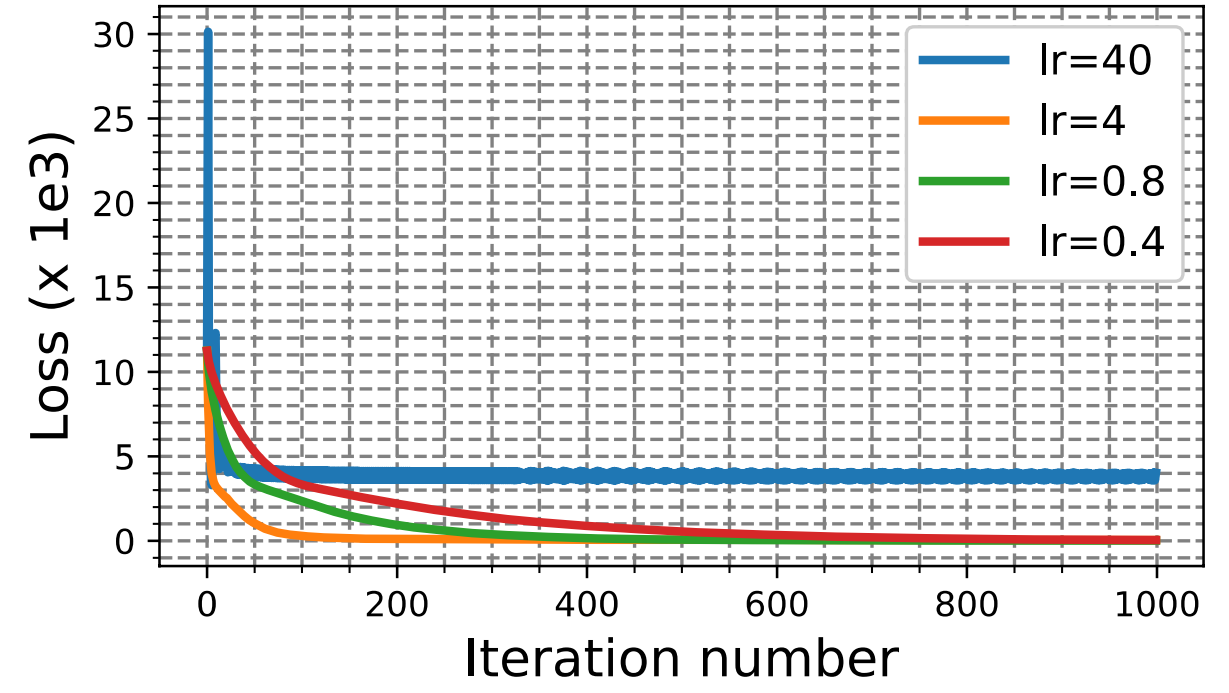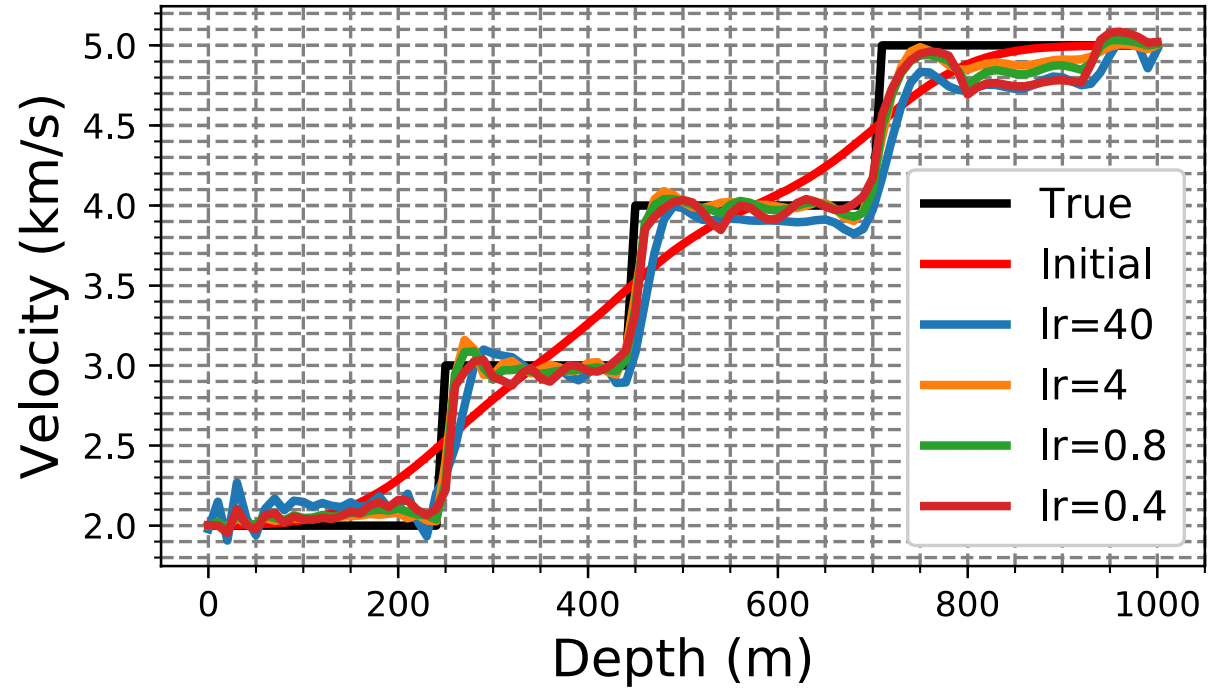
$\alpha \in (10,100)$

$$\boldsymbol{r}_k = \beta \cdot \boldsymbol{r}_{k-1} + (1 - \beta) \cdot \boldsymbol{g}_k^2$$

$$\tilde{\boldsymbol{r}}_k = \boldsymbol{r}_k / (1 - \beta^k)$$

$$\boldsymbol{v}_k = \boldsymbol{v}_{k-1} - \boxed{\frac{\alpha}{\sqrt{\tilde{\boldsymbol{r}}_k} + \epsilon}} \cdot \boldsymbol{g}_k$$

$$\alpha \in (1,10)$$

$$\boldsymbol{m}_k = \beta_1 \cdot \boldsymbol{m}_{k-1} + (1 - \beta_1) \cdot \boldsymbol{g}_k$$

$$\boldsymbol{r}_k = \beta_2 \cdot \boldsymbol{r}_{k-1} + (1 - \beta_2) \cdot \boldsymbol{g}_k^2$$

$$\tilde{\boldsymbol{m}}_k = \boldsymbol{m}_k / (1 - \beta_1^k)$$

$$\tilde{\boldsymbol{r}}_k = \boldsymbol{r}_k / (1 - \beta_2^k)$$

$$\boldsymbol{v}_k = \boldsymbol{v}_{k-1} - \boxed{\frac{\alpha}{\sqrt{\tilde{\boldsymbol{r}}_k} + \epsilon}} \cdot \tilde{\boldsymbol{m}}_k$$

$$\alpha \in (10, 100)$$

$$\boldsymbol{m}_k = \beta_1 \cdot \boldsymbol{m}_{k-1} + (1 - \beta_1) \cdot \boldsymbol{g}_k$$

$$\boldsymbol{r}_k = \beta_2 \cdot \boldsymbol{r}_{k-1} + (1 - \beta_2) \cdot \boldsymbol{g}_k^2$$

$$\tilde{\boldsymbol{m}}_k = \boldsymbol{m}_k / (1 - \beta_1^k)$$

$$\tilde{\boldsymbol{r}}_k = \boldsymbol{r}_k / (1 - \beta_2^k)$$

$$\boldsymbol{v}_k = \boldsymbol{v}_{k-1} - \boxed{\frac{\alpha}{\sqrt{\tilde{\boldsymbol{r}}_k} + \epsilon}} \cdot \tilde{\boldsymbol{m}}_k$$

$$\alpha \in (10, 100)$$

# Synthetic test on Marmousi

- Marmousi: 11 shots (12Hz)

❑ Conclusions:
- ▪ We illustrated a self-designed RNN framework for forward and inverse seismic modelling.
- ▪ Proved that FWI is a special case of a machine learning process.
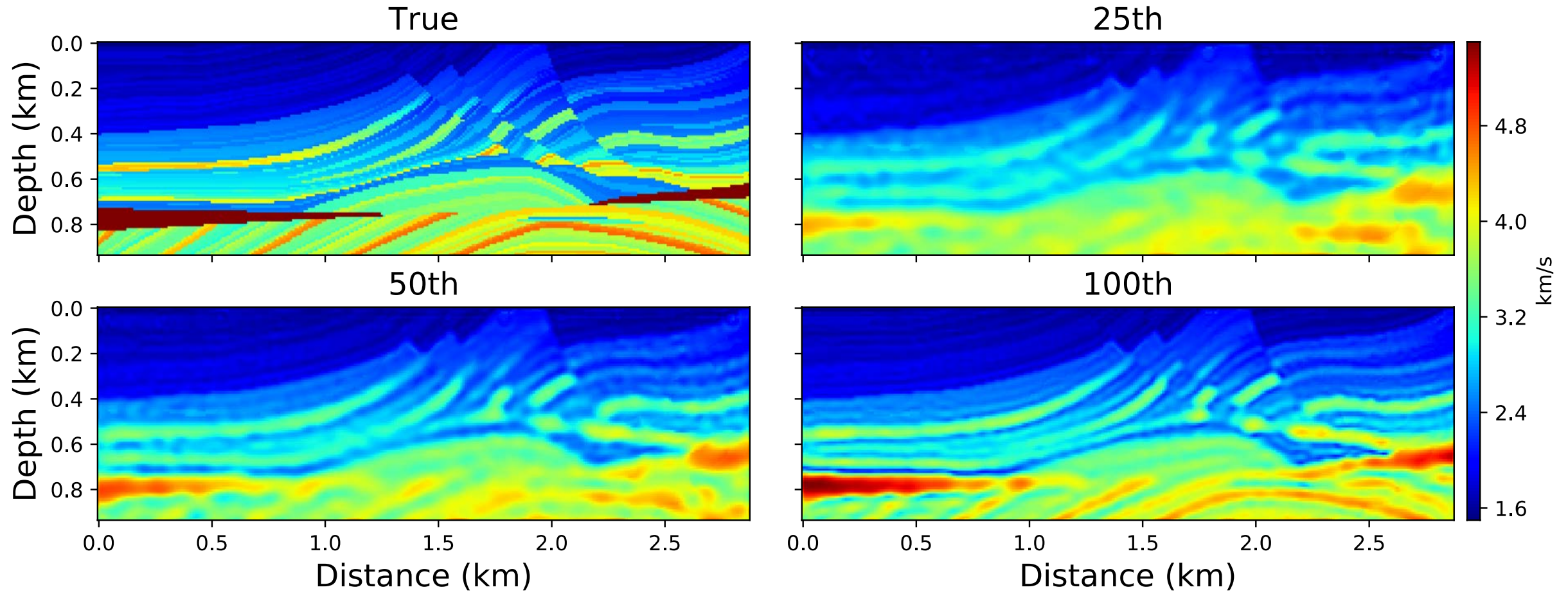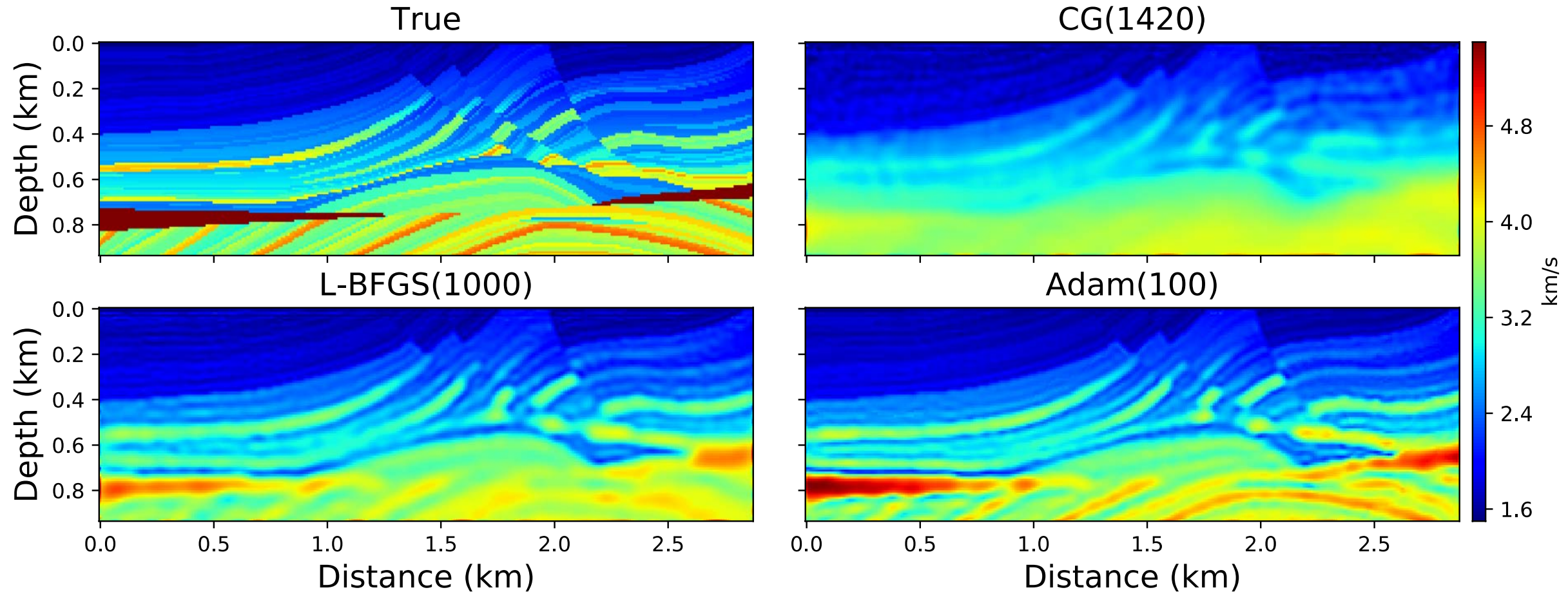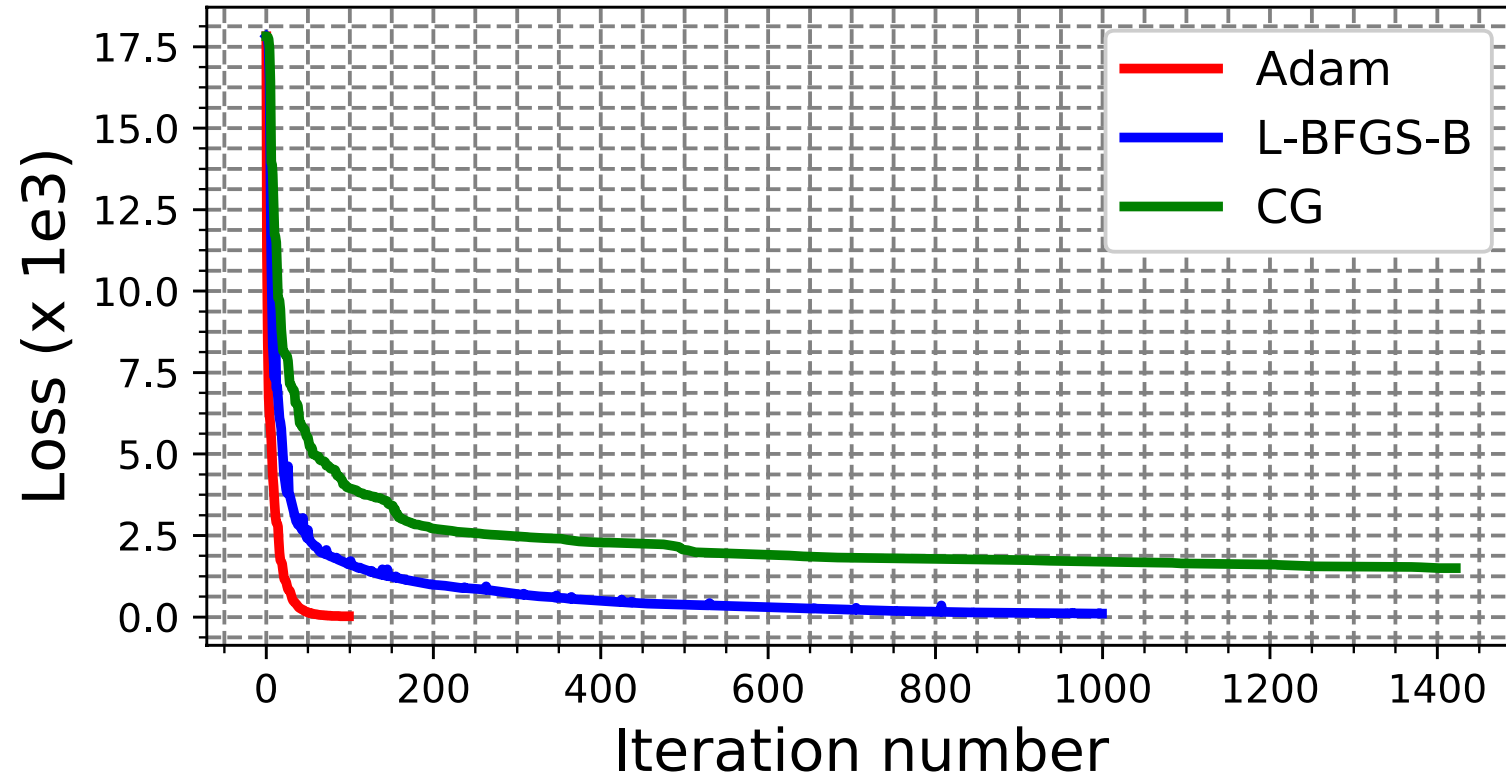- ▪ Best learning rate ranges of gradient-based algorithms were analyzed and investigated for velocity model building.
- ▪ The efficiency of gradient-based and non-linear optimization algorithms are compared and discussed.

❑ Future works:
- ▪ A theory-guided neural network.
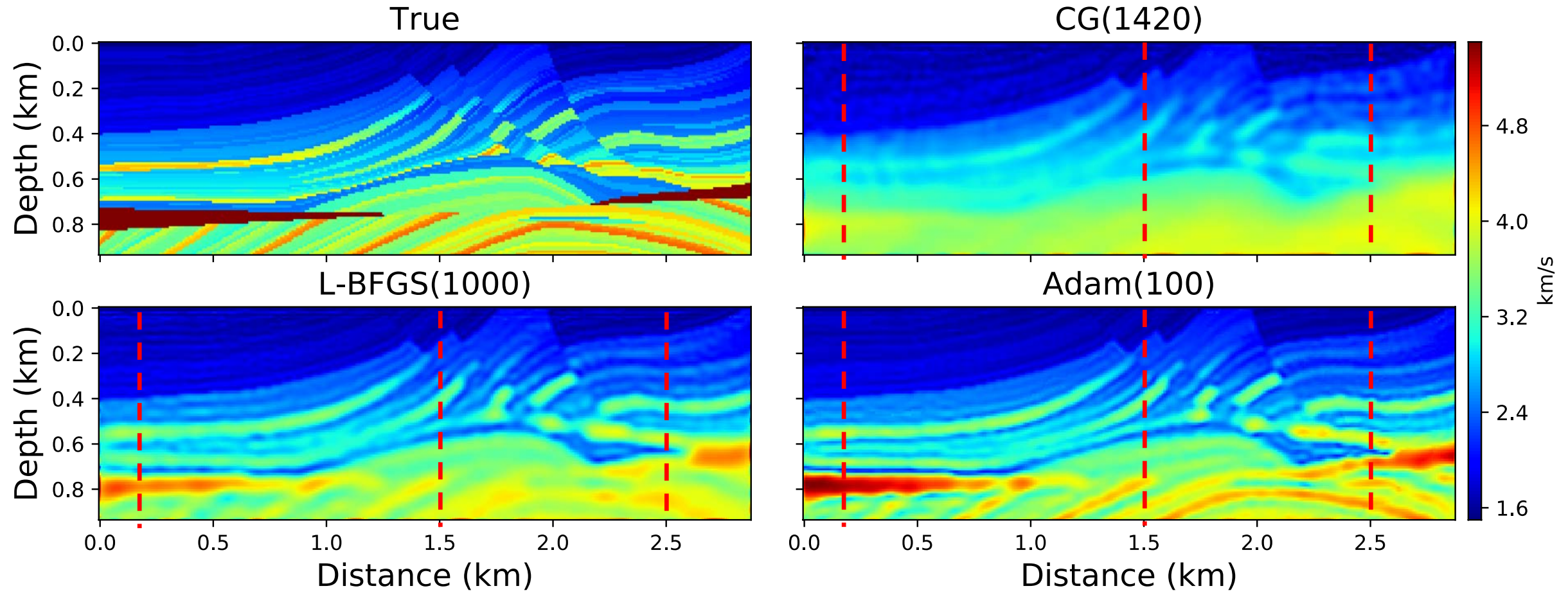- ▪ A physical-teaching training process.

# Acknowledgements

- All CREWES sponsors
- NSERC (CRDPJ 46117913)
- Kiki Xu, Sergio Romahn, Lei Yang, Xin Fu
- All CREWES staff and students

Comments & Questions ?