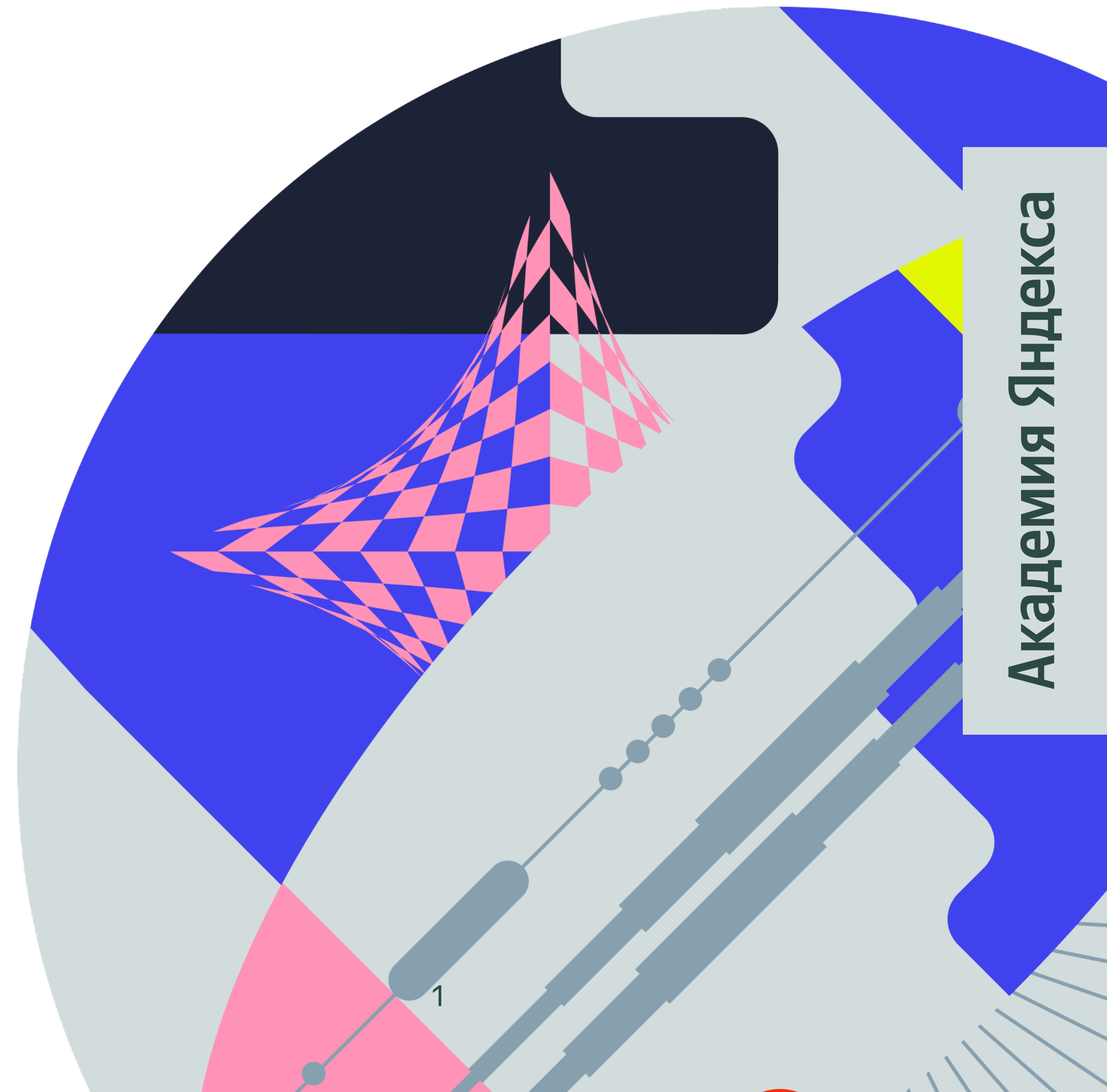


# Visual Generative Modeling 2025

**Nikita Starodubcev**

*Yandex Research*



# Offtop

DM <sup>?</sup> == FM

# Offtop

In the FM, we have the following:

$$\boldsymbol{x} = (1 - t)\boldsymbol{x}_0 + \boldsymbol{\epsilon}t, \quad t \in [0,1], \quad d\boldsymbol{x} = \boldsymbol{u}_\theta(\boldsymbol{x}, t)dt, \quad \boldsymbol{u}_\theta := \boldsymbol{\epsilon}_\theta - \boldsymbol{x}_\theta$$

Let's show that FM is a special case of DMs. The basic equation in DMs:

$$d\boldsymbol{x} = \left[ \boldsymbol{x}f(t) - \frac{1}{2}g(t)^2 \nabla_{\boldsymbol{x}} \log p_t(\boldsymbol{x}) \right] dt.$$

Which  $f(t)$  and  $g(t)$  should be to obtain the FM forward process? To this end, let's remember the connection between  $f(t)$ ,  $g(t)$  and statistics of the forward process (expectation and variance).

$$d\mathbb{E}_{\boldsymbol{x}}\boldsymbol{x} = f(t)\mathbb{E}_{\boldsymbol{x}}\boldsymbol{x}dt, \quad \mathbb{E}_{\boldsymbol{x}}\boldsymbol{x}(0) = \boldsymbol{x}_0.$$

$$d\mathbb{D}_{\boldsymbol{x}}\boldsymbol{x} = \left( 2f(t)\mathbb{D}_{\boldsymbol{x}}\boldsymbol{x} + g(t)^2 \right) dt, \quad \mathbb{D}_{\boldsymbol{x}}\boldsymbol{x}(0) = 0.$$

# Offtop

$$f(t) = -\frac{1}{1-t} . \qquad g(t)^2 = \frac{2t}{1-t} .$$

$$\nabla_{\boldsymbol{x}} \log p_t(\boldsymbol{x}) \approx -\frac{1}{t^2} \left( \boldsymbol{x} - (1-t)\boldsymbol{x}_\theta \right) = -\frac{\boldsymbol{\epsilon}_\theta}{t}$$

$$\mathrm{d}\boldsymbol{x} = \left[ -\boldsymbol{x} \frac{1}{1-t} + \frac{1}{2} \frac{2t}{1-t} \frac{\boldsymbol{\epsilon}_\theta}{t} \right] \mathrm{d}t,$$

$$\mathrm{d}\boldsymbol{x} = \frac{1}{1-t} \left[ -(1-t)\boldsymbol{x}_\theta - \boldsymbol{\epsilon}_\theta t + \boldsymbol{\epsilon}_\theta \right] \mathrm{d}t,$$

$$\mathrm{d}\boldsymbol{x} = \boldsymbol{u}_\theta \mathrm{d}t$$

# Offtop

DM == FM

FM is special case of DM with specific noising process

# Lecture 3 | Numerical solvers

1. Summary (recap from the previous lectures)
2. Basics of numerical solution of ODE
3. DPM-solver
4. Sampling schedules

Noising process  $\longrightarrow$  ODE (SDE)  $\longrightarrow$  Training of NN  $\longrightarrow$  Inference using solver

|              | Noising process  | ODE   | Parameterization  |
|--------------|--|---|---|
| General case | $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \mathbf{z},$ $t \in [0, T], \mathbf{z} \sim \mathcal{N}(0, 1)$                  | $d\mathbf{x} = \left[ \mathbf{x} f(t) - \frac{1}{2} g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt$ $f(t) = \frac{d \log \alpha_t}{dt}, g^2(t) = \frac{d \sigma_t^2}{dt} - 2 \frac{d \log \alpha_t}{dt} \sigma_t^2$ | $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = - \frac{\mathbf{x} - \alpha_t \mathbb{E} \mathbf{x}_0   \mathbf{x}}{\sigma_t^2}$ $\boldsymbol{\epsilon} = \frac{\mathbf{x} - \alpha_t \mathbb{E} \mathbf{x}_0   \mathbf{x}}{\sigma_t}, \quad \mathbb{E} \mathbf{x}_0   \mathbf{x} \approx \mathbf{x}_0$ |
| VP           | $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \mathbf{z}$ $t \in [0, T], \mathbf{z} \sim \mathcal{N}(0, 1)$ | $d\mathbf{x} = \frac{1}{2\alpha_t} [\mathbf{x} + \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] d\alpha_t$   | $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = - \frac{\mathbf{x} - \sqrt{\alpha_t} \mathbb{E} \mathbf{x}_0   \mathbf{x}}{1 - \alpha_t}$   |
| VE           | $\mathbf{x}_t = \mathbf{x}_0 + t \mathbf{z}$ $t \in [0, T], \mathbf{z} \sim \mathcal{N}(0, 1)$                                   | $d\mathbf{x} = -t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) dt$  | $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = - \frac{\mathbf{x} - \mathbb{E} \mathbf{x}_0   \mathbf{x}}{t^2}$  |
| FM<br>(RF)   | $\mathbf{x}_t = (1 - t) \mathbf{x}_0 + t \mathbf{z}$ $t \in [0, 1], \mathbf{z} \sim \mathcal{N}(0, 1)$                           | $d\mathbf{x} = - \frac{1}{1 - t} [\mathbf{x} + t \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt$ $d\mathbf{x} = \mathbf{v}(\mathbf{x}, t) dt$   | $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = - \frac{\mathbf{x} - (1 - t) \mathbb{E} \mathbf{x}_0   \mathbf{x}}{t^2}$ $\mathbf{v} = \boldsymbol{\epsilon} - \mathbb{E} \mathbf{x}_0   \mathbf{x}$  |

$$1. \quad \mathbb{E}_{\mathbf{x}_0, t, \mathbf{x}_t} \|\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|$$

$$2. \quad \mathbb{E}_{\mathbf{x}_0, t, \mathbf{x}_t} \|s_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|$$

$$3. \quad \mathbb{E}_{\mathbf{x}_0, t, \mathbf{x}_t} \|\mathbf{x}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|$$

# 1. Conditional expectation

$$\mathbb{E}_{\mathbf{x}_0, t, \mathbf{x}_t} \|\mathbf{x}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|$$

$$\mathbf{x}_\theta(\mathbf{x}, t) \approx \mathbb{E} \mathbf{x}_0 | \mathbf{x} = \int \mathbf{x}_0 p_t(\mathbf{x}_0 | \mathbf{x}) d\mathbf{x}_0.$$

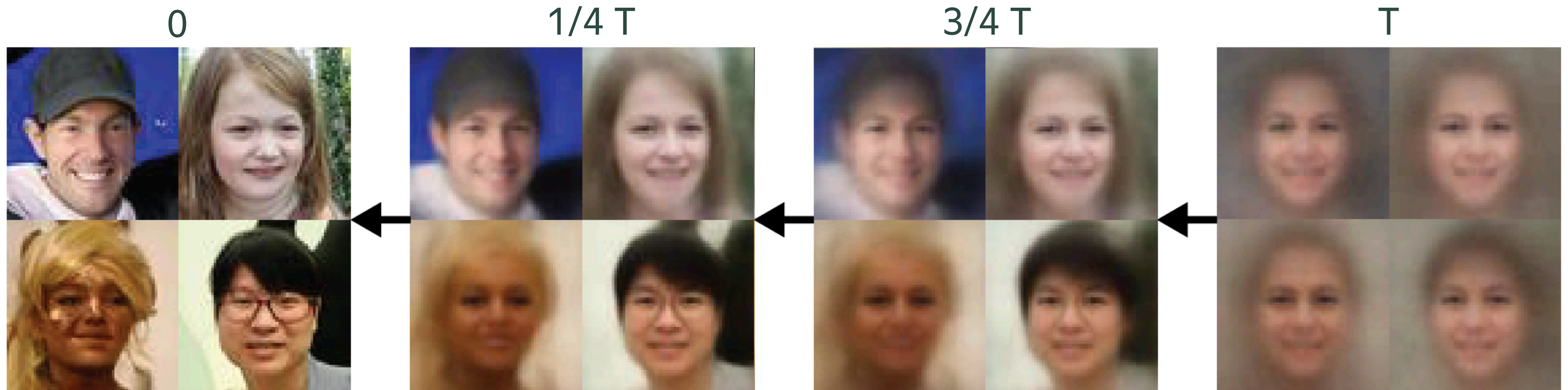
$$p_t(\mathbf{x}_0 | \mathbf{x}) = \frac{p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0)}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0}, \quad p_{data}(\mathbf{x}_0) = \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{x}_0 - \mathbf{x}_0^j).$$

$$\mathbb{E} \mathbf{x}_0 | \mathbf{x} = \int \mathbf{x}_0 \frac{p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0)}{\int p_t(\mathbf{x} | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0} d\mathbf{x}_0 = \dots = \frac{\sum_{j=1}^N \mathbf{x}_0^j \mathcal{N}(\mathbf{x} | \alpha_t \mathbf{x}_0^j, \sigma_t^2)}{\sum_{j=1}^N \mathcal{N}(\mathbf{x} | \alpha_t \mathbf{x}_0^j, \sigma_t^2)}.$$



# 1. Conditional expectation

$$\mathbb{E} \mathbf{x}_0 | \mathbf{x} = \frac{\sum_{j=1}^N \mathbf{x}_0^j \mathcal{N}(\mathbf{x} | \alpha_t \mathbf{x}_0^j, \sigma_t^2)}{\sum_{j=1}^N \mathcal{N}(\mathbf{x} | \alpha_t \mathbf{x}_0^j, \sigma_t^2)}.$$



We make the images less and less averaged, coarse-to-fine generation

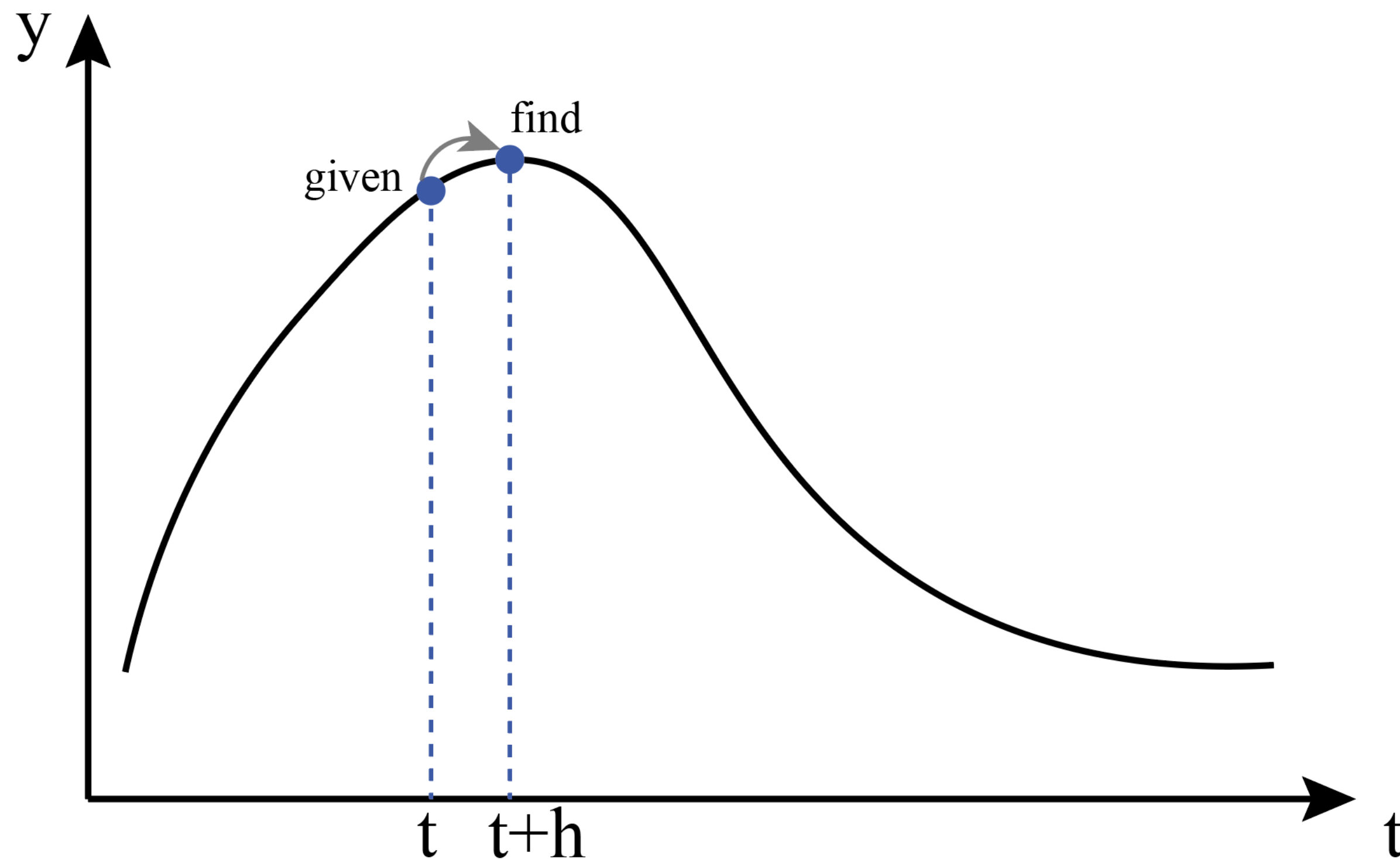
# 2. Basics of numerical solution of ODE

Cauchy problem

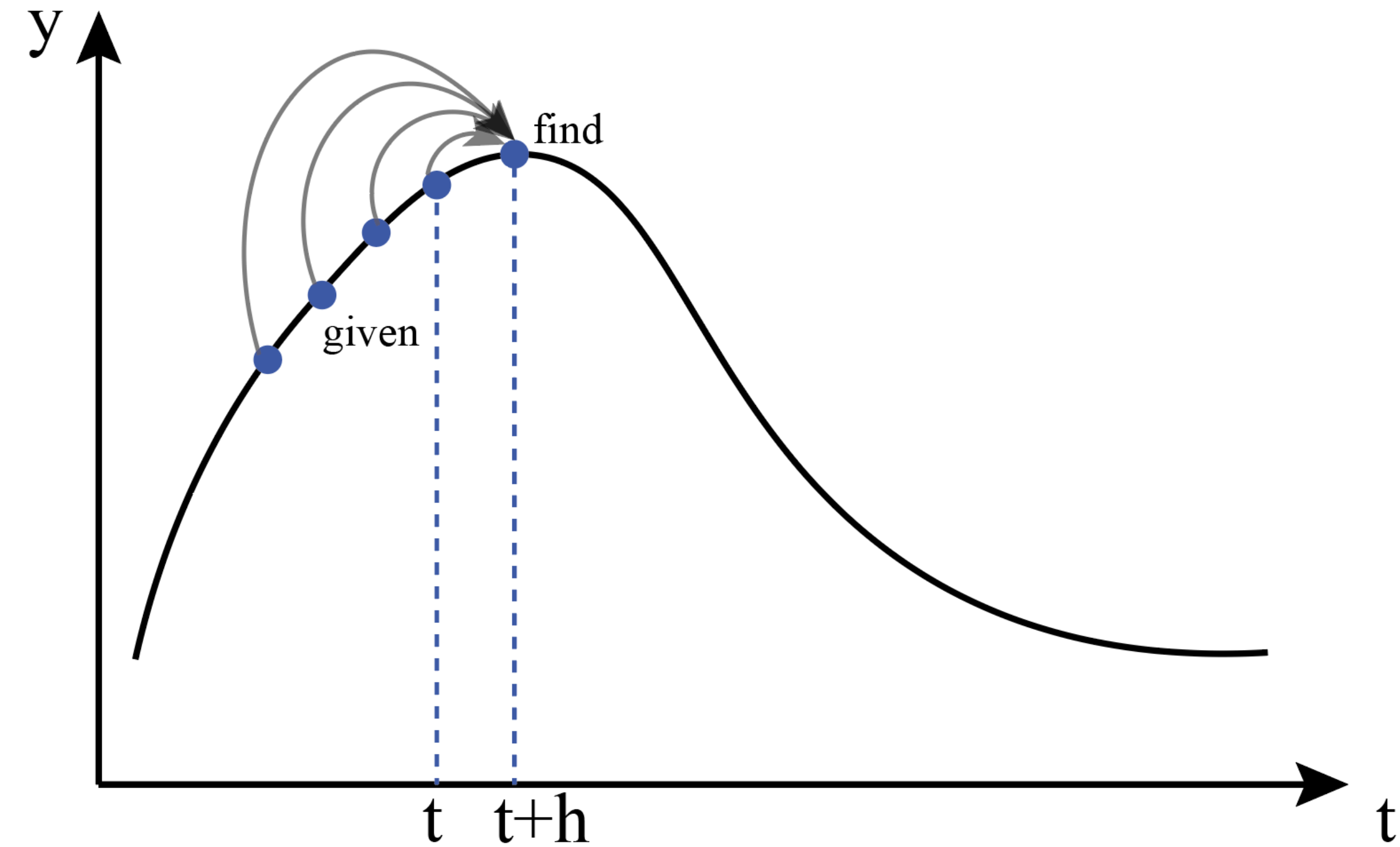
$$dy = f(t, y(t))dt, y(0) = y_0$$

Numerical solvers  $[t_0, \dots, t_N], h = t_{i+1} - t_i$

Singlestep  $y_{i+1} = y_i + \Phi(y_i, h)$



Multistep  $y_{i+1} = y_i + \Phi(y_i, y_{i-1}, \dots, h)$



Singlestep  $y_{i+1} = y_i + \Phi(y_i, h)$

Multistep  $y_{i+1} = y_i + \Phi(y_i, y_{i-1}, \dots, h)$

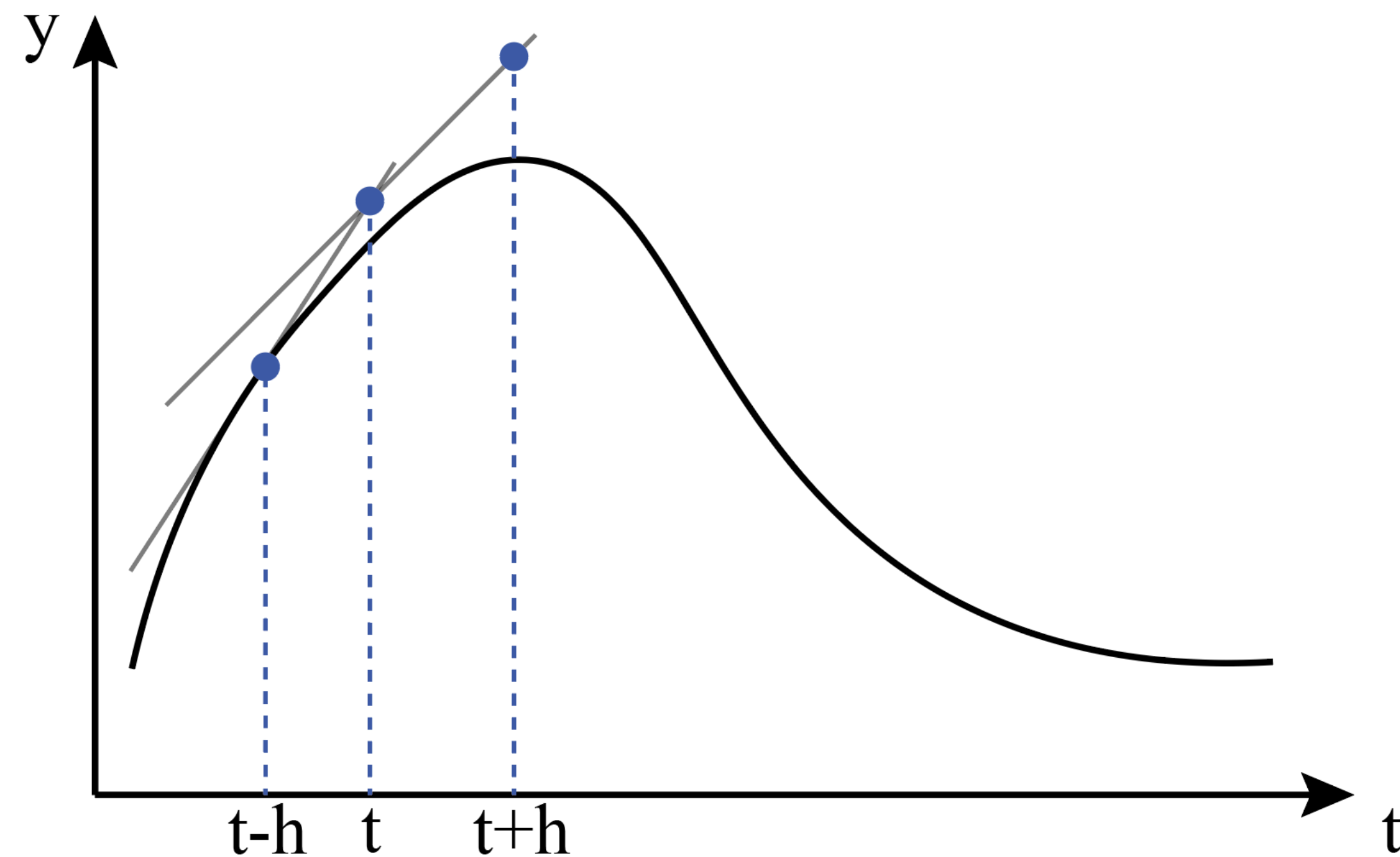
Implicit  $y_{i+1} = y_i + \Phi(y_{i+1}, y_i, h)$

Implicit  $y_{i+1} = y_i + \Phi(y_{i+1}, y_i, y_{i-1}, \dots, h)$

Examples:

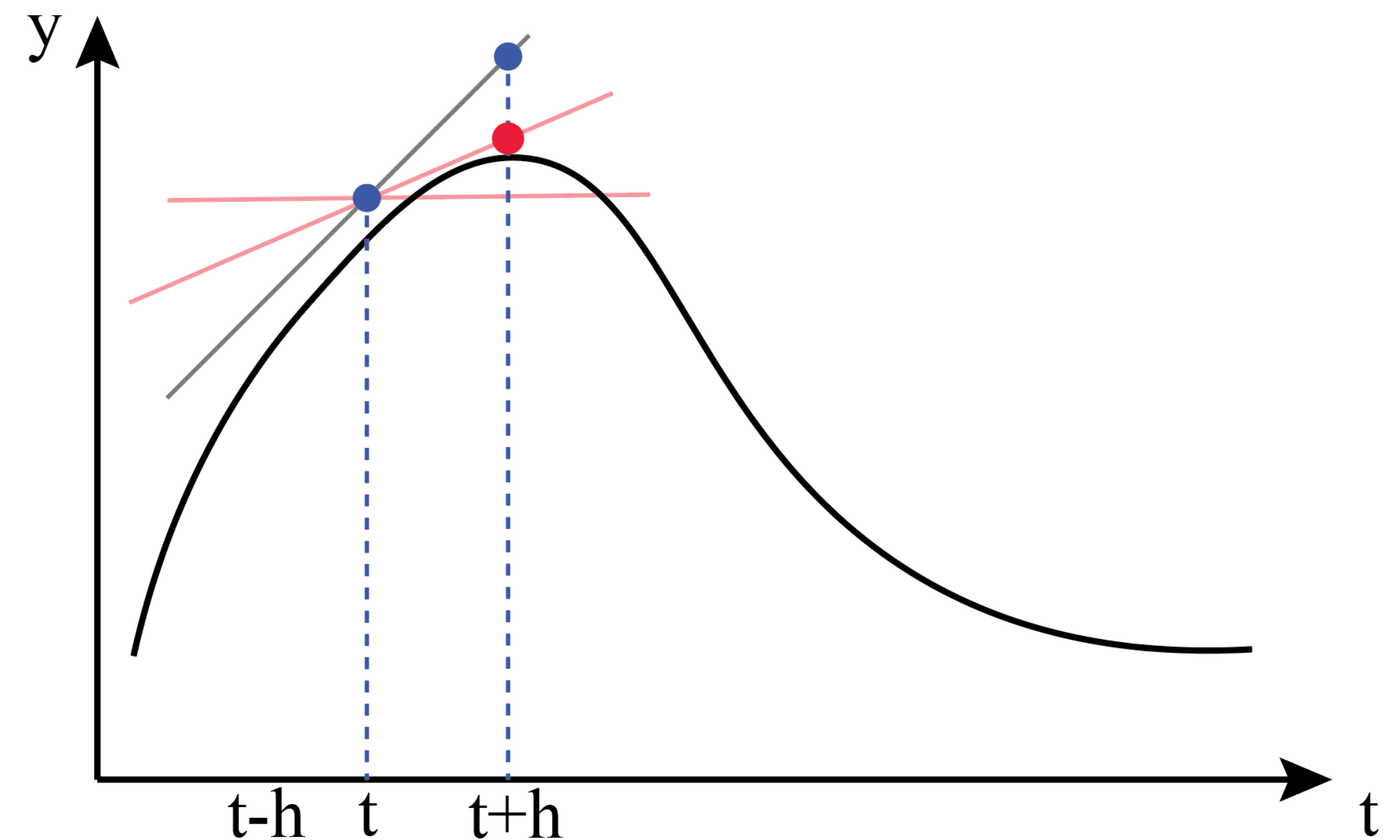
\* Euler (singlestep, explicit)

$$y_{i+1} = y_i + hf(y_i, t_i)$$



\* Trapezoidal rule (singlestep, implicit)

$$y_{i+1} = y_i + \frac{h}{2} (f(y_i, t_i) + f(y_{i+1}, t_{i+1}))$$



How to understand how accurate is your solver?

1. Local truncation error (we assume previous steps are exact)

$$\tau_{i+1} = \|y_{i+1} - y(t_{i+1})\|$$

\* Euler

$$y(t_{i+1}) = y_i + hy'_i + \frac{h^2}{2}y''_i + O(h^3)$$

$$y_{i+1} = y_i + hf(y_i, t_i), f(y_i, t_i) = y'_i$$

$$\tau_{i+1} = y_i + hy'_i - \left( y_i + hy'_i + \frac{h^2}{2}y''_i + O(h^3) \right) = -\frac{h^2}{2}y''_i + O(h^3) \sim O(h^2)$$

\* Trapezoidal rule

$$\tau_{i+1} \sim O(h^3)$$

## 2. Global truncation error (accumulative error after N steps)

$$e_N = \|y_N - y(t_N)\| \quad e_N \sim N\tau \sim \frac{\tau}{h}$$

\* Euler

$$e_N \sim O(h)$$

$$\int_{t_i}^{t_{i+1}} dy = \int_{t_i}^{t_{i+1}} f(t, y(t)) dt$$

$$y_{i+1} = y_i + \int_{t_i}^{t_{i+1}} f(t, y(t)) dt$$

$$f(t, y(t)) = \sum_{k=0} \frac{h^k}{k!} f^{(k)}(t_i, y_i) = f(t_i, y_i) + O(h)$$

$$y_{i+1} = y_i + f(t_i, y_i) + O(h^2)$$

\* Trapezoidal rule

$$e_N \sim O(h^2)$$

Solver converges if:

1. LTE  $\rightarrow 0$ , as  $h \rightarrow 0$

2. It is stable

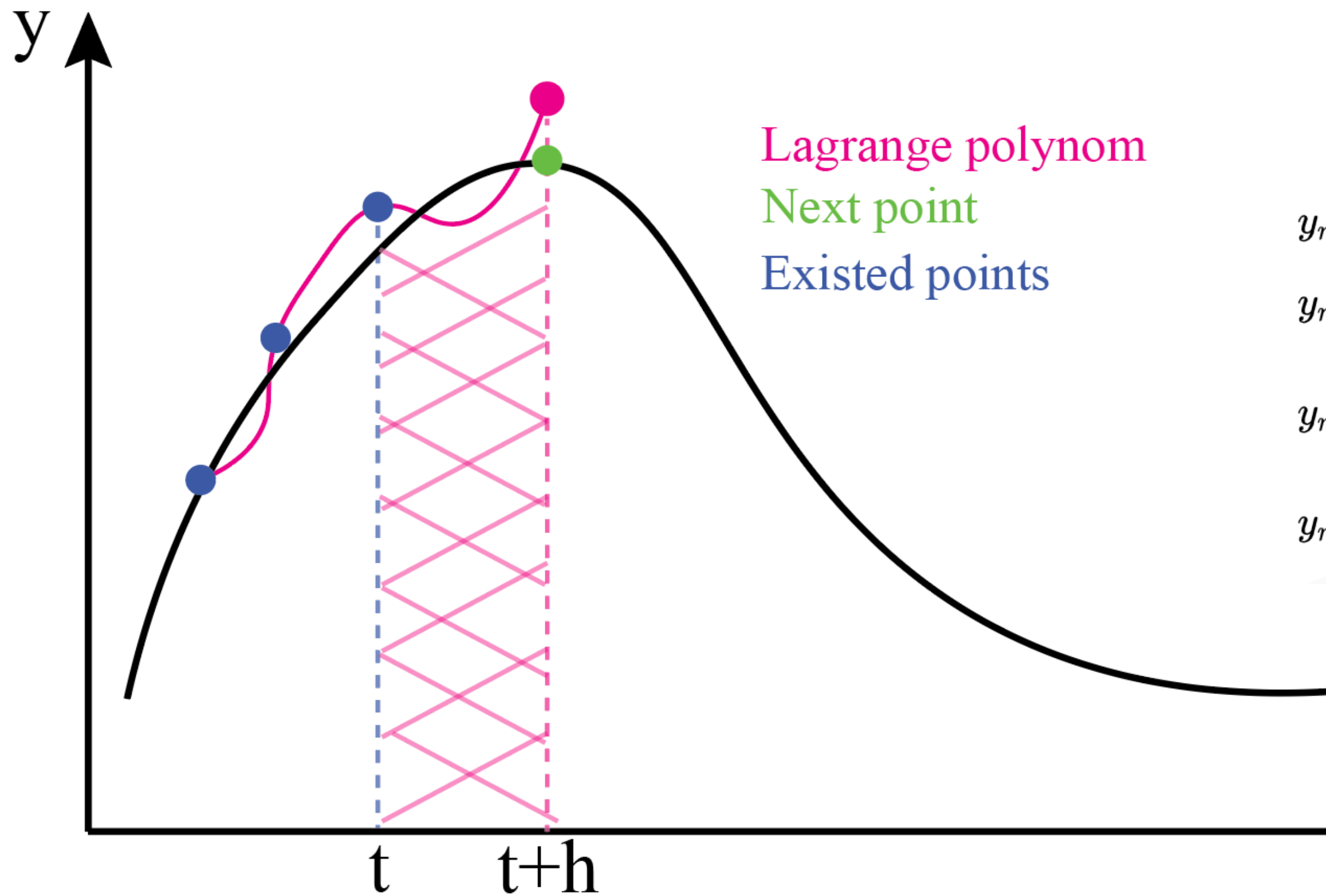
small perturbations do not grow exponentially.

# Multistep explicit methods (Adams-Bashforth methods)

Use past solution values and past derivative evaluations to compute the next step

$$y_{i+1} = y_i + \int_{t_i}^{t_{i+1}} f(t, y(t)) dt \quad \text{Cannot integrate } f(t, y(t))$$

$$f(t, y(t)) \approx p_{s-1}(t), \quad p_{s-1}(t) = \sum_{j=0}^{s-1} \frac{(-1)^{s-j-1} f(t_{n+j}, y_{n+j})}{j!(s-j-1)!h^{s-1}} \prod_{i=0}^{s-1} (t - t_{n+i}) \quad \text{Lagrange polynomial}$$



$$y_{n+1} = y_n + hf(t_n, y_n), \quad (\text{This is the Euler method})$$

$$y_{n+2} = y_{n+1} + h \left( \frac{3}{2} f(t_{n+1}, y_{n+1}) - \frac{1}{2} f(t_n, y_n) \right),$$

$$y_{n+3} = y_{n+2} + h \left( \frac{23}{12} f(t_{n+2}, y_{n+2}) - \frac{16}{12} f(t_{n+1}, y_{n+1}) + \frac{5}{12} f(t_n, y_n) \right),$$

$$y_{n+4} = y_{n+3} + h \left( \frac{55}{24} f(t_{n+3}, y_{n+3}) - \frac{59}{24} f(t_{n+2}, y_{n+2}) + \frac{37}{24} f(t_{n+1}, y_{n+1}) - \frac{9}{24} f(t_n, y_n) \right),$$

## Multistep vs singlestep (DPM-solver specific)

$$y_{i+1} = y_i + \int_{t_i}^{t_{i+1}} f(t, y(t)) dt$$

$$f(t, y(t)) = \sum_{k=0} \frac{h^k}{k!} f^{(k)}(t_i, y_i) = f(t_i, y_i) + hf'(t_i, y_i) + O(h^2)$$

How to approximate  $f'(t_i, y_i)$  ?

Singlestep (use next intermediate point)

$$f'(t_i, y_i) \approx \frac{f(t_{i+\delta_i}, y_{i+\delta_i}) - f(t_i, y_i)}{t_{i+\delta_i} - t_i}$$

Multistep (use previous point)

$$f'(t_i, y_i) \approx \frac{f(t_i, y_i) - f(t_{i-1}, y_{i-1})}{h}$$



# Summary

- I) Singlestep explicit  $y_{i+1} = y_i + \Phi(y_i, h)$   
\* Euler (1 NFE per step, 1st order)  
\* RK 2 (2 NFE per step, 2nd order)  
Use the intermediate steps
- II) Singlestep implicit  $y_{i+1} = y_i + \Phi(y_{i+1}, y_i, h)$   
\* Backward Euler  
\* RK 2 implicit, Heun (2 NFE per step, 2nd order)  
Use the next steps
- III) Multistep explicit  $y_{i+1} = y_i + \Phi(y_i, y_{i-1}, \dots, h)$   
\* Adams-Bashforth (1 NFE, any order)  
Use the previous steps

- IV) Multistep implicit  $y_{i+1} = y_i + \Phi(y_{i+1}, y_i, y_{i-1}, h)$   
\* Adams-Moulton  
Use the next and previous steps
- V) Multistep + Singlestep  
\* General linear methods  
Use the previous and intermediate steps



## 2. DPM-solver (VP)

$$d\mathbf{x} = \frac{1}{2\alpha_t} \left[ \mathbf{x} + \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] d\alpha_t$$

DDIM (Denoising Diffusion Implicit Models)

$$\mathbf{x}_t = \sqrt{\alpha_t} \left( \frac{\mathbf{x}_s - \sqrt{1 - \alpha_s} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_s, s)}{\sqrt{\alpha_s}} \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_s, s)$$

What is DDIM? Just Euler?

## 2. DPM-solver (VP)

$$d\mathbf{x} = \frac{1}{2\alpha_t} [\mathbf{x} + \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] d\alpha_t$$

DDIM (Denoising Diffusion Implicit Models)

$$\mathbf{x}_t = \sqrt{\alpha_t} \left( \frac{\mathbf{x}_s - \sqrt{1 - \alpha_s} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_s, s)}{\sqrt{\alpha_s}} \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_s, s)$$

What is DDIM? Just Euler?

$$d\mathbf{x} = \frac{1}{2\alpha_t} \left[ \mathbf{x} - \frac{\boldsymbol{\epsilon}(\mathbf{x}, t)}{\sqrt{1 - \alpha_t}} \right] d\alpha_t$$

$$\mathbf{x}_t = \frac{\alpha_t + \alpha_s}{2\alpha_s} \mathbf{x}_s - \frac{\alpha_t - \alpha_s}{2\alpha_s \sqrt{1 - \alpha_s}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_s, s)$$

$$\mathbf{x}_t = \mathbf{x}_s + \int_s^t \frac{1}{2\alpha_{\tau}} \left[ \mathbf{x} - \frac{\boldsymbol{\epsilon}(\mathbf{x}, \tau)}{\sqrt{1 - \alpha_t}} \right] d\alpha_{\tau}$$

$$d\mathbf{x} = \frac{1}{2\alpha_t} \left[ \mathbf{x} - \frac{\epsilon(\mathbf{x}, t)}{\sqrt{1 - \alpha_t}} \right] d\alpha_t$$

$$d\mathbf{x} = \left[ \frac{\mathbf{x}}{2\alpha_t} - \frac{\epsilon(\mathbf{x}, t)}{2\alpha_t\sqrt{1 - \alpha_t}} \right] d\alpha_t, s \rightarrow t$$

$$\phi(t, s)d\mathbf{x} = \phi(t, s) \left[ \frac{\mathbf{x}}{2\alpha_t} - \frac{\epsilon(\mathbf{x}, t)}{2\alpha_t\sqrt{1 - \alpha_t}} \right] d\alpha_t,$$

$$d [\mathbf{x}\phi(t, s)] = \phi(t, s) \left[ \frac{\mathbf{x}}{2\alpha_t} - \frac{\epsilon(\mathbf{x}, t)}{2\alpha_t\sqrt{1 - \alpha_t}} \right] d\alpha_t,$$

$$d\phi(t, s) = \phi(t, s) \frac{1}{2\alpha_t} d\alpha_t$$

$$\int_s^t \frac{d\phi}{\phi} = \int_s^t \frac{d\alpha}{2\alpha} \rightarrow \ln \phi(t, s) = \frac{1}{2} (\ln \alpha_t - \ln \alpha_s)$$

$$\phi(t, s) = \sqrt{\frac{\alpha_t}{\alpha_s}}$$

$$\mathbf{x}_t = \sqrt{\frac{\alpha_t}{\alpha_s}} \mathbf{x}_s - \frac{1}{2} \int_s^t \phi(t, \tau) \frac{\epsilon(\mathbf{x}, \tau)}{\alpha_\tau \sqrt{1 - \alpha_\tau}} d\alpha_\tau,$$

$$\mathbf{x}_t = \sqrt{\frac{\alpha_t}{\alpha_s}} \mathbf{x}_s - \frac{\sqrt{\alpha_t}}{2} \int_s^t \frac{\epsilon(\mathbf{x}, \tau)}{\sqrt{\alpha_\tau} \alpha_\tau \sqrt{1 - \alpha_\tau}} d\alpha_\tau,$$

$$\int_s^t \frac{\epsilon(\mathbf{x}, \tau)}{\sqrt{\alpha_\tau} \alpha_\tau \sqrt{1 - \alpha_\tau}} d\alpha_\tau \approx \epsilon(\mathbf{x}_s, s) \int_s^t \frac{1}{\sqrt{\alpha_\tau} \alpha_\tau \sqrt{1 - \alpha_\tau}} d\alpha_\tau$$

$$\int \frac{d\alpha}{\sqrt{\alpha}\alpha\sqrt{1-\alpha}}$$

$$1-\alpha=\cos^2u,\quad \alpha=\sin^2u,\quad d\alpha=2\sin u\cos udu$$

$$\int \frac{2du}{\sin^2u}=-2\operatorname{ctg}u=-2\operatorname{ctg}\arcsin\sqrt{\alpha}=-2\frac{\sqrt{1-\alpha}}{\sqrt{\alpha}}$$

DDIM = DPM-1

$$\boldsymbol{x}_t=\sqrt{\frac{\alpha_t}{\alpha_s}}\boldsymbol{x}_s-\frac{\sqrt{\alpha_t}}{2}\left(-2\frac{\sqrt{1-\alpha_t}}{\sqrt{\alpha_t}}+2\frac{\sqrt{1-\alpha_s}}{\sqrt{\alpha_s}}\right)\boldsymbol{\epsilon}(\boldsymbol{x}_s,s)$$

$$\boldsymbol{x}_t=\sqrt{\alpha_t}\left(\frac{\boldsymbol{x}_s-\sqrt{1-\alpha_s}\boldsymbol{\epsilon}_{\theta}(\boldsymbol{x}_s,s)}{\sqrt{\alpha_s}}\right)+\sqrt{1-\alpha_t}\boldsymbol{\epsilon}(\boldsymbol{x}_s,s)$$

$$\boldsymbol{x}_t = \sqrt{\frac{\alpha_t}{\alpha_s}} \boldsymbol{x}_s - \frac{\sqrt{\alpha_t}}{2} \int_s^t \frac{\boldsymbol{\epsilon}(\boldsymbol{x}, \tau)}{\sqrt{\alpha_\tau} \alpha_\tau \sqrt{1 - \alpha_\tau}} d\alpha_\tau,$$

$$\boldsymbol{\epsilon}(\boldsymbol{x}, \tau) = \sum_{k=0}^{n-1} \frac{(\alpha_\tau - \alpha_s)^k}{k!} \boldsymbol{\epsilon}^{(k)}(\boldsymbol{x}_s, s)$$

Singlestep (use the intermediate point)

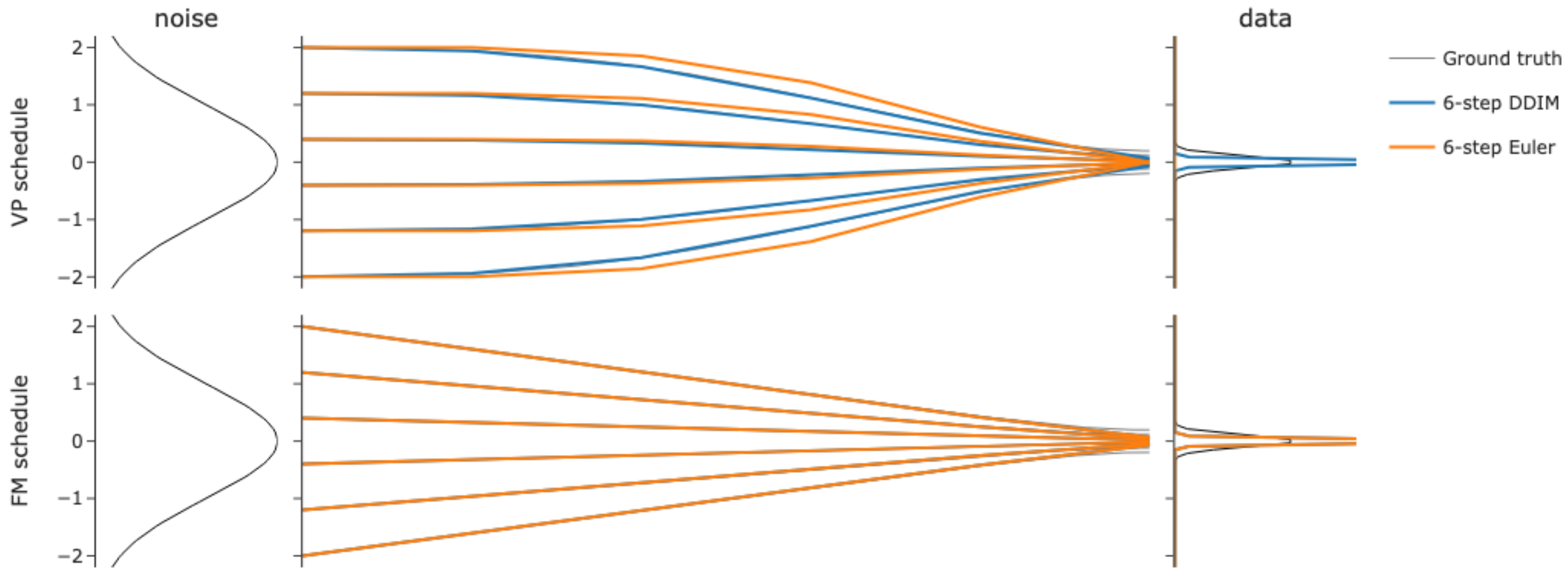
$$\boldsymbol{\epsilon}^{(1)}(\boldsymbol{x}_s, s) \approx \frac{\boldsymbol{\epsilon}(\boldsymbol{x}_{s-\delta_s}, s - \delta_s) - \boldsymbol{\epsilon}(\boldsymbol{x}_s, s)}{\alpha_{s-\delta_s} - \alpha_s}$$

Multistep (use the previous point)

$$\boldsymbol{\epsilon}^{(1)}(\boldsymbol{x}_s, s) \approx \frac{\boldsymbol{\epsilon}(\boldsymbol{x}_s, s) - \boldsymbol{\epsilon}(\boldsymbol{x}_{s+h}, s + h)}{\alpha_s - \alpha_{s+h}}$$

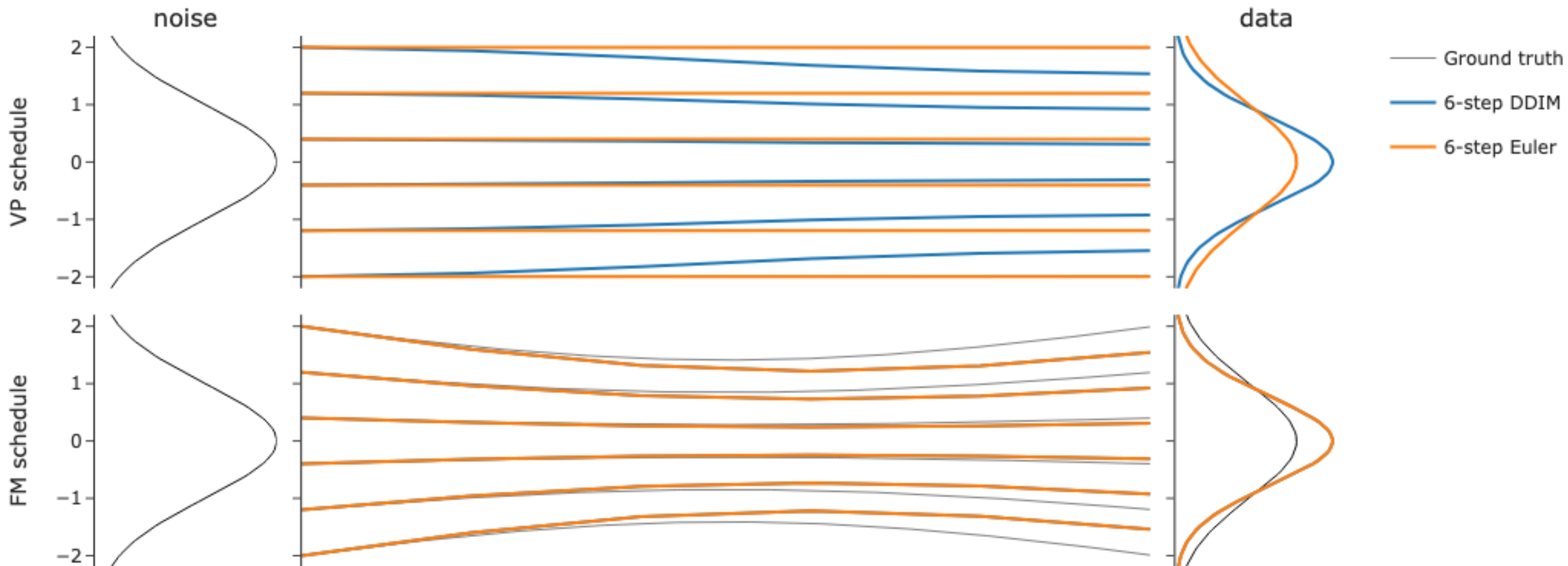
# 2. Sampling schedules

Variance Preserving vs Flow Matching schedules for varying data distributions



# 2. Sampling schedules

Variance Preserving vs Flow Matching schedules for varying data distributions





## 2. Sampling schedules

