



TOBIAS BUCK

27.9.2024 TOBIAS.BUCK@IWR.UNI-HEIDELBERG.DE

SCIENTIFIC ML

OUTLINE OF TODAY

1. Scientific ML
 - A. Neural ODEs
 - B. Symbolic Regression
 - C. Operator Learning
 - D. Physics Informed Neural Networks
 - E. Hamiltonian Neural Networks

WHERE TO FIND THE LECTURE MATERIAL



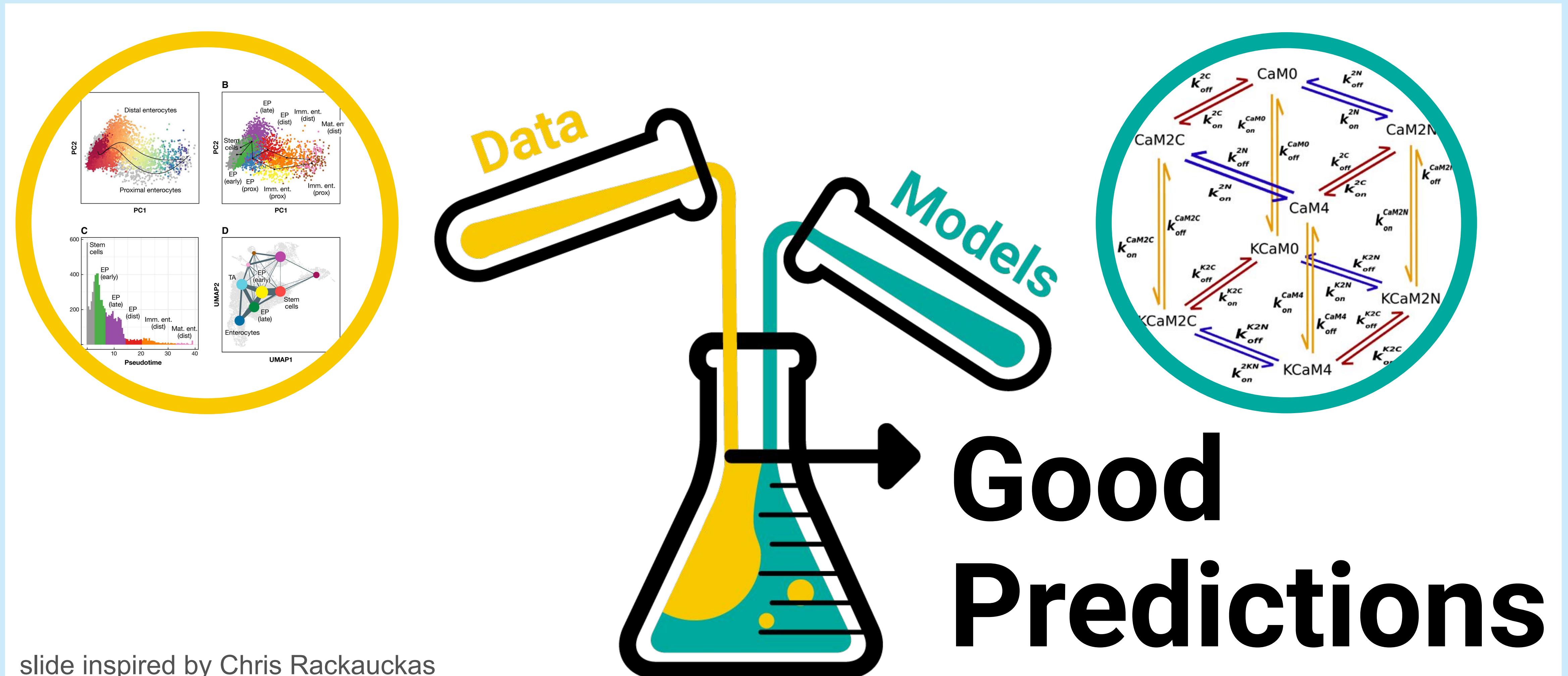
<https://github.com/TobiBu/graddays>

SCIENTIFIC MACHINE LEARNING

»Scientific machine learning is the burgeoning field combining techniques of machine learning into traditional scientific computing and mechanistic modeling.« [e. g. Zubov et al, 2021]

SCIENTIFIC MACHINE LEARNING IS MODEL-BASED DATA-EFFICIENT ML

How do we simultaneously use both sources of knowledge?



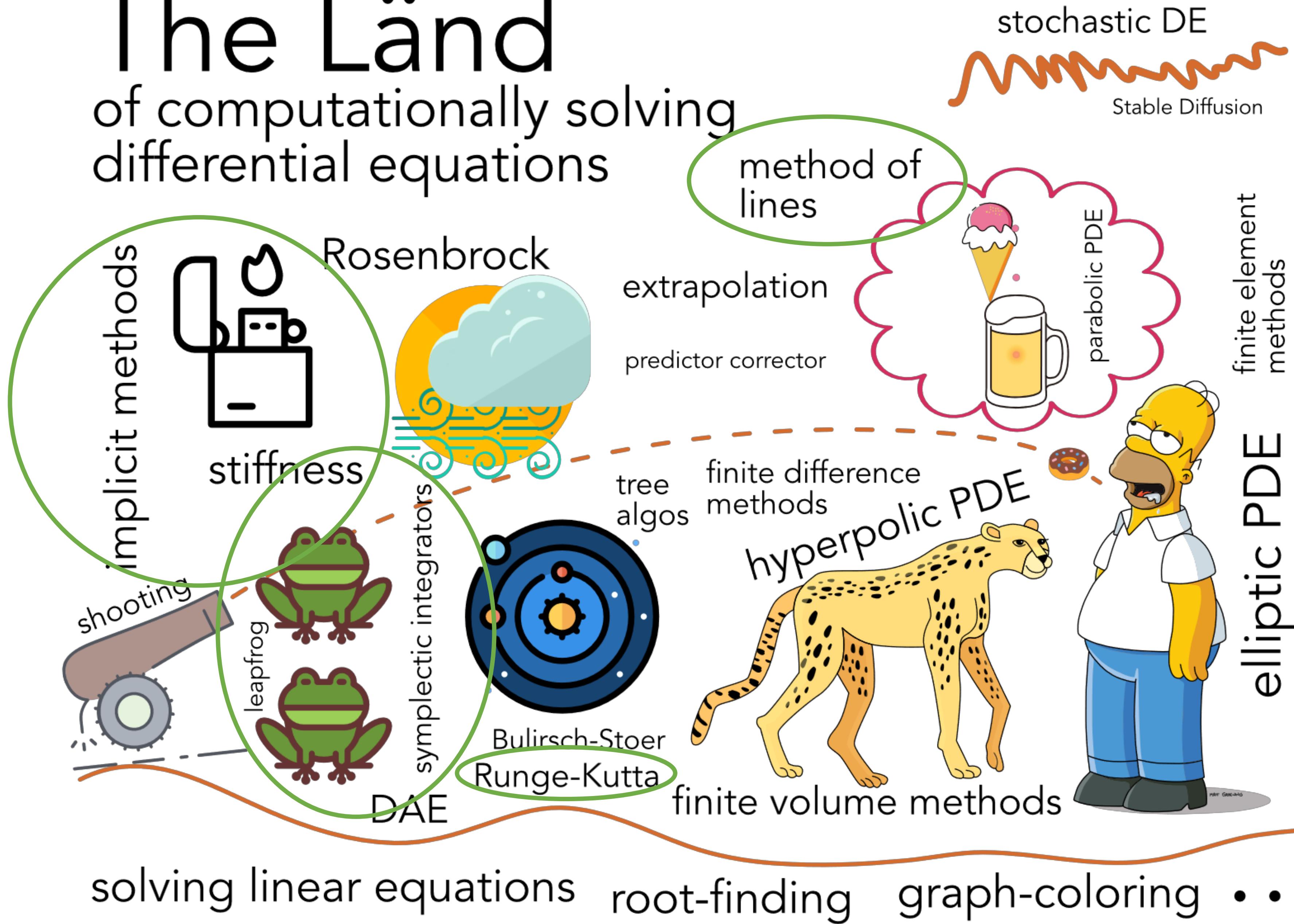
**Core Point of SciML:
the more structure that is encoded into the
model, the less there is to learn and the more
data efficient the algorithm is!**

Traditional methods for solving DEs

it's bio

The Länd

of computationally solving
differential equations



THE SOLUTION OF ODES

Explicit Euler Method:

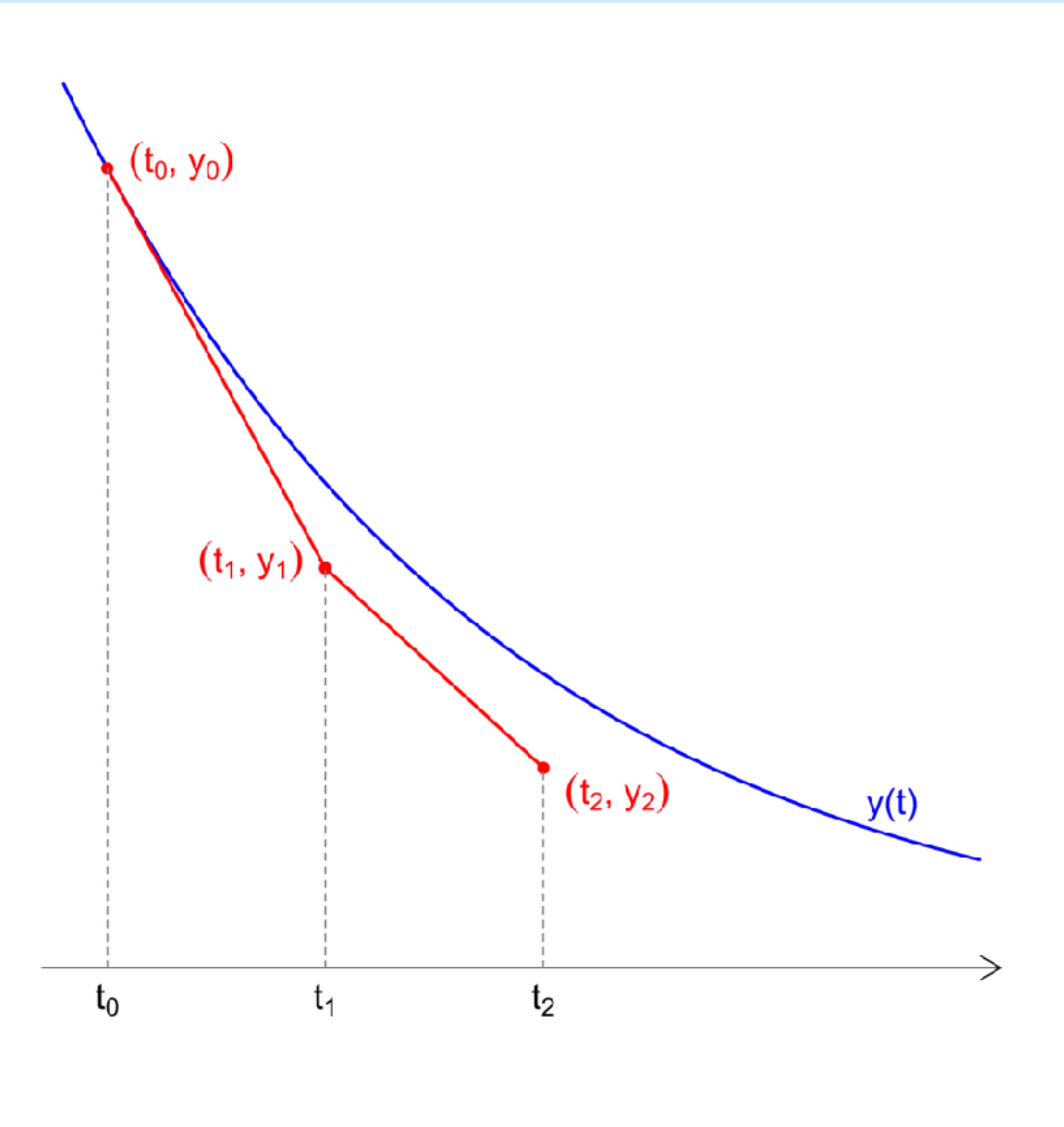
$$y^{n+1} = y^n + \frac{\partial y}{\partial t} \Delta t$$

Order of the solver:

Taylor expansion

$$y^{n+1} = y^n + \frac{\partial y^n}{\partial t} \Delta t + \mathcal{O}(\Delta t^2)$$

→ first order accurate

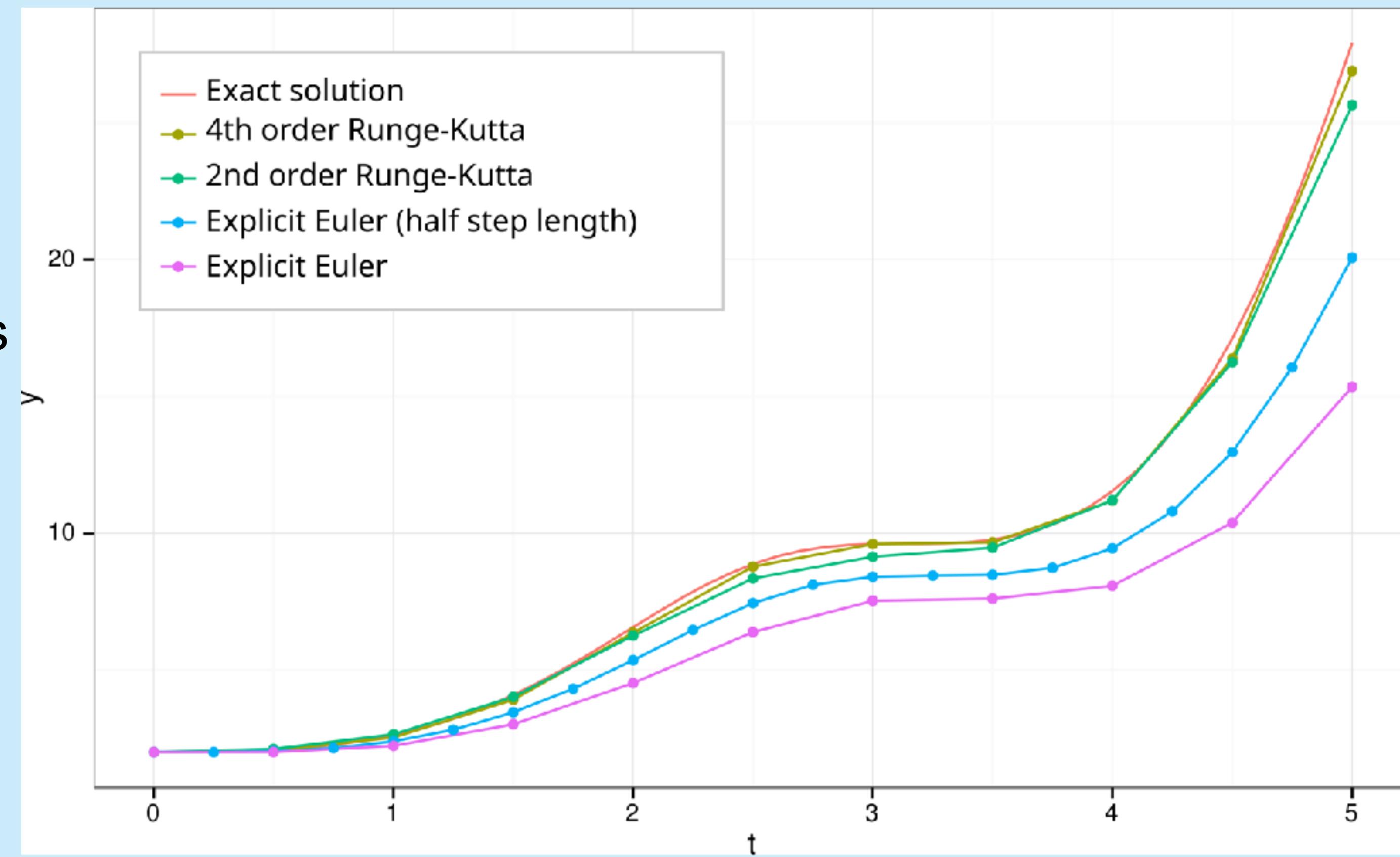


HIGHER ORDER SOLVERS

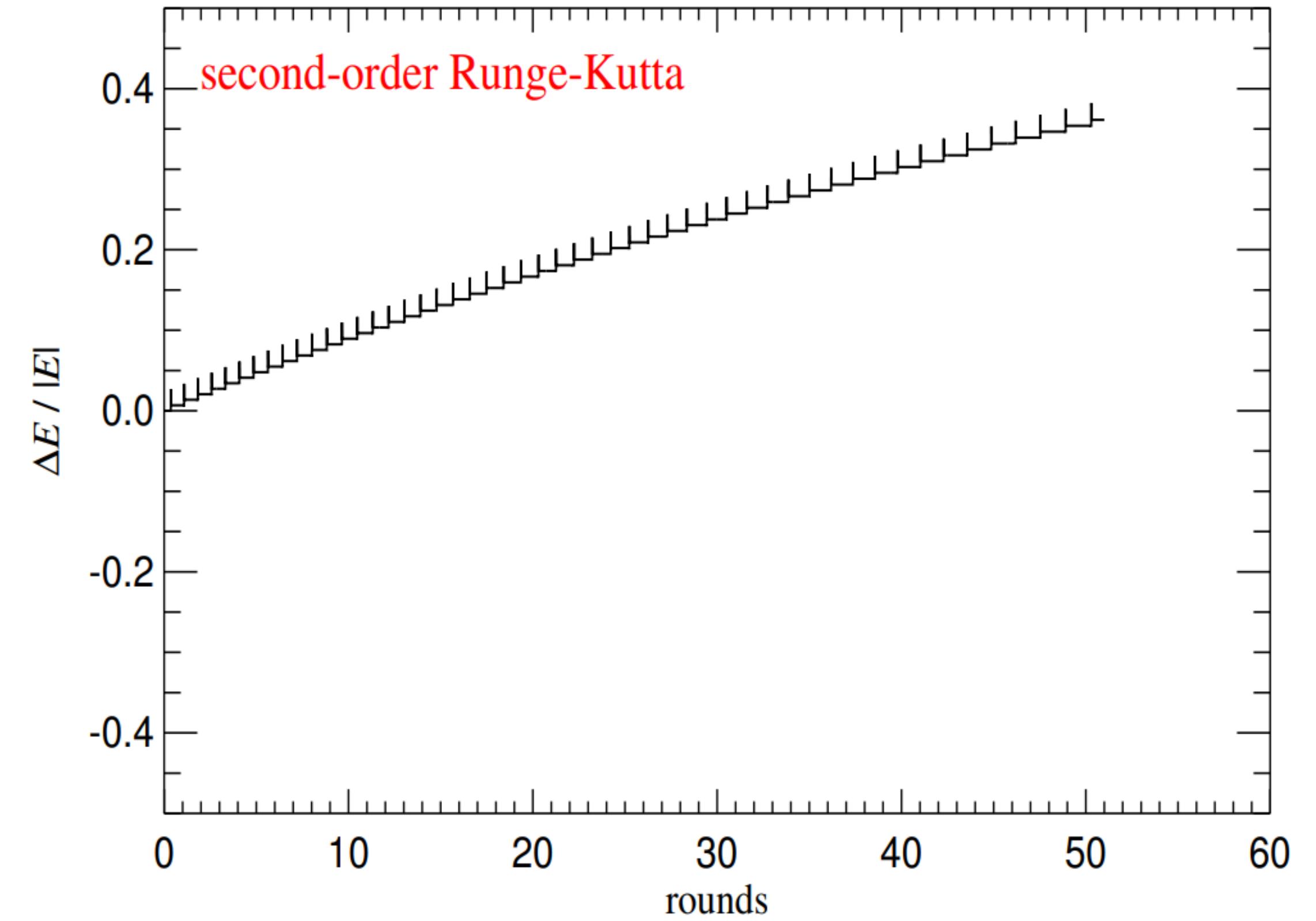
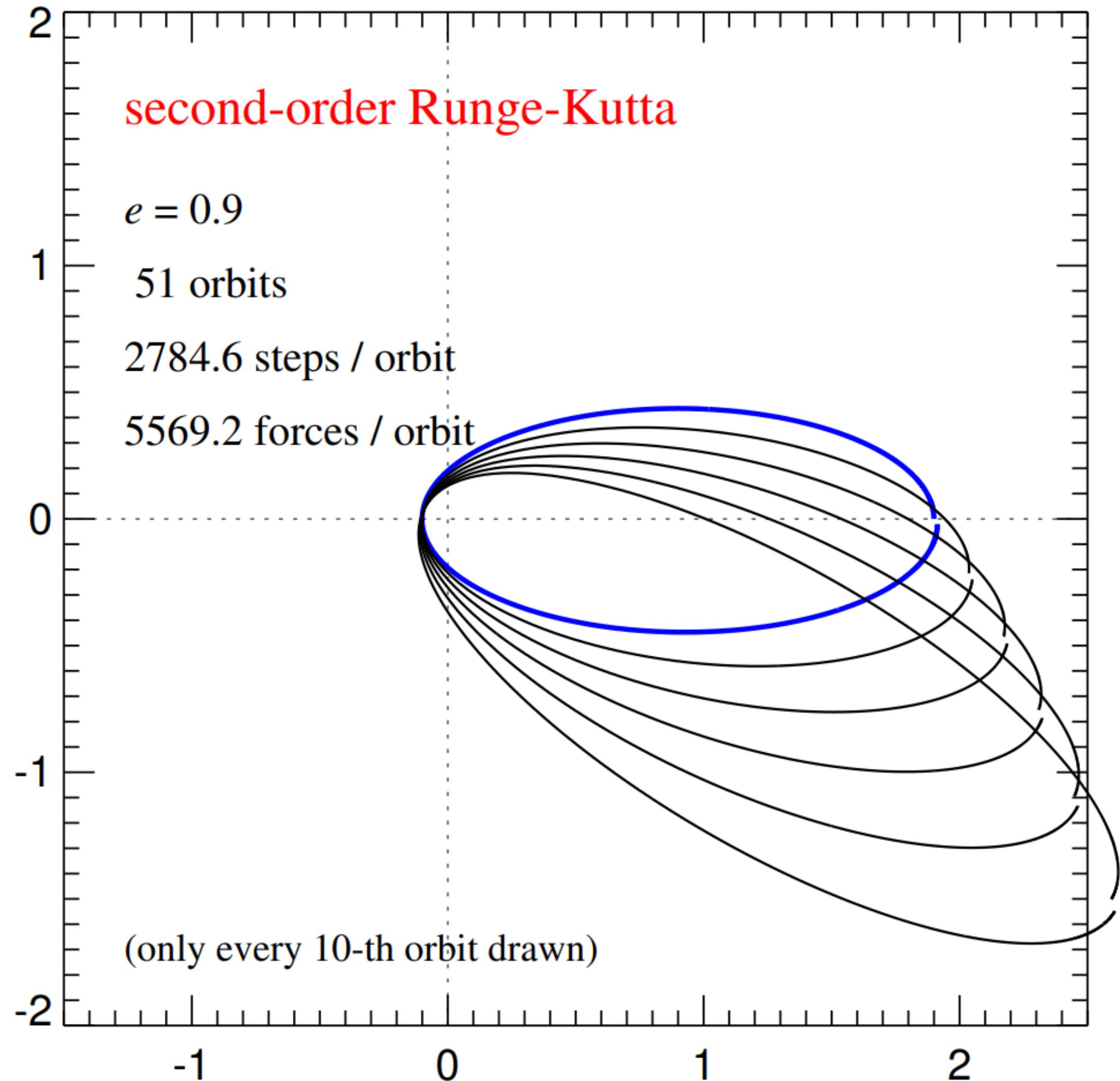
Runge-Kutta-Method:
weighted combination of simple derivatives

$$y^{n+1} = y^n + h \sum_{i=1}^s b_i k_i$$

$$\begin{aligned}k_1 &= f(t_n, y_n), \\k_2 &= f(t_n + c_2 h, y_n + (a_{21} k_1) h), \\k_3 &= f(t_n + c_3 h, y_n + (a_{31} k_1 + a_{32} k_2) h), \\&\vdots \\k_s &= f(t_n + c_s h, y_n + (a_{s1} k_1 + a_{s2} k_2 + \cdots + a_{s,s-1} k_{s-1}) h).\end{aligned}$$



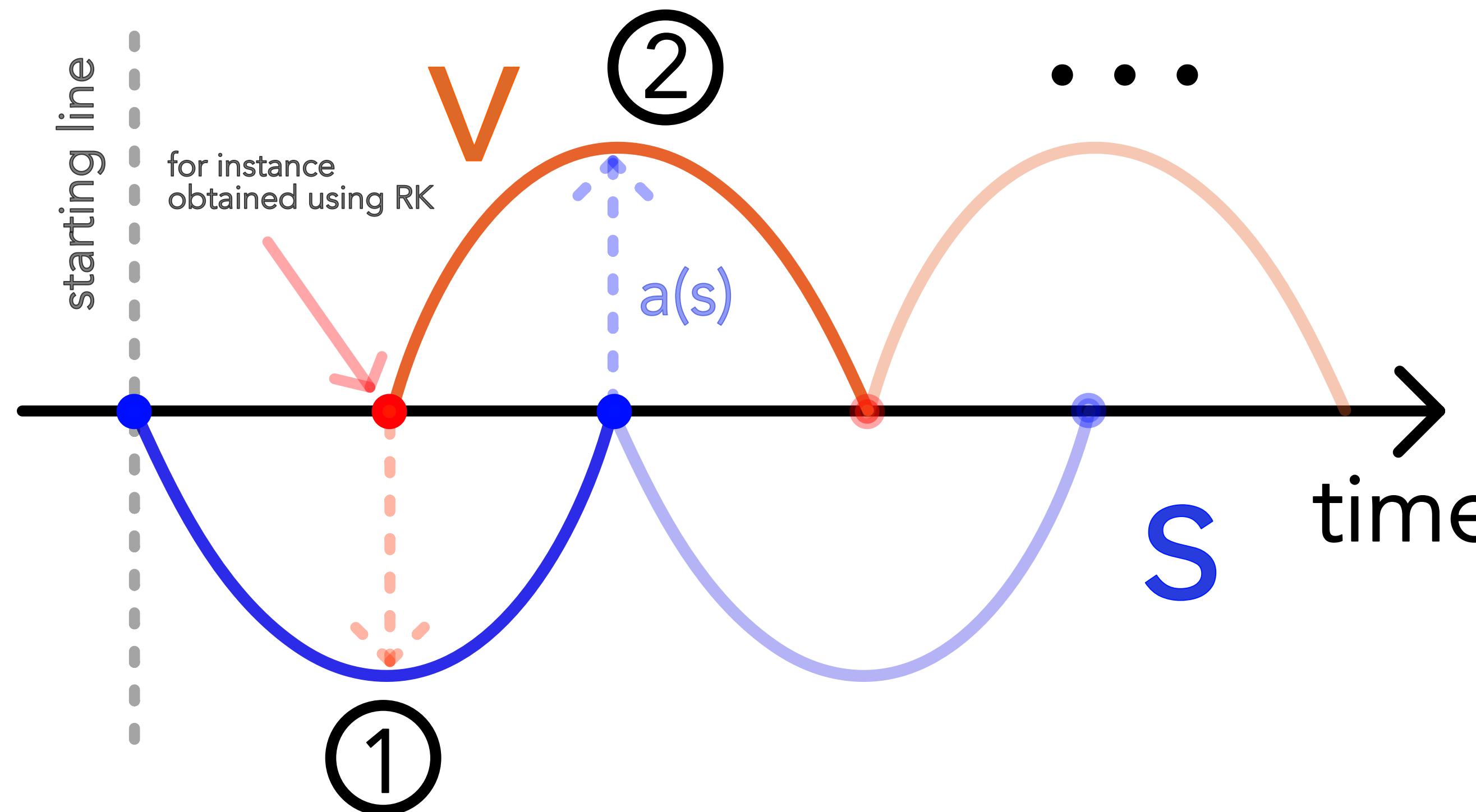
SOLUTION OF THE TWO-BODY PROBLEM



Energy error in each step!

SOLUTION OF THE TWO-BODY PROBLEM

Leapfrog



Symplectic integrator with good energy conservation properties; time reversibility; exact angular momentum conservation

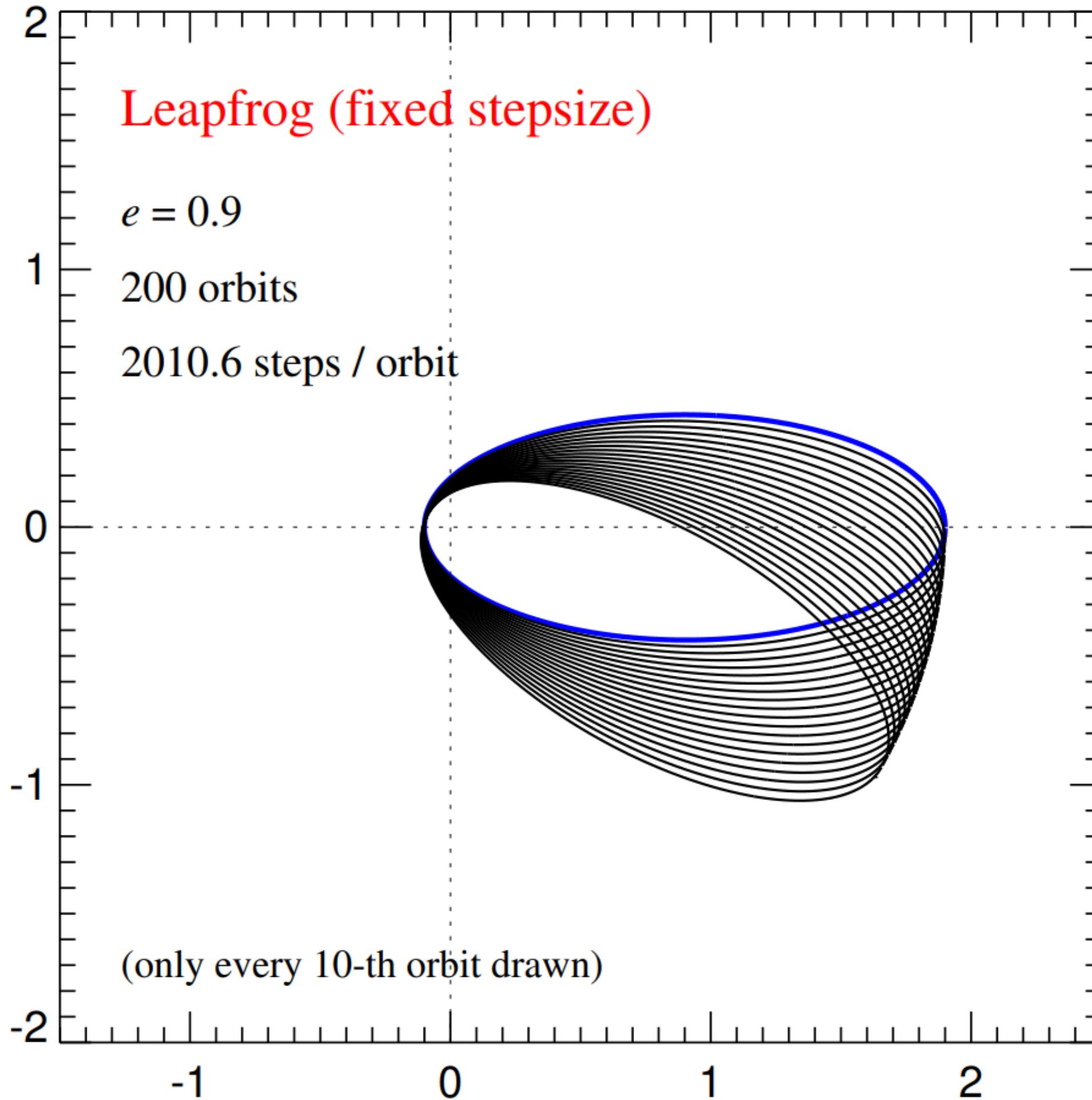
1. update location:

$$\begin{aligned}\underline{s}\left(t + \frac{1}{2}\Delta t\right) = \\ \underline{s}\left(t - \frac{1}{2}\Delta t\right) + \underline{v}(t)\Delta t + \\ \mathcal{O}(\Delta t^3)\end{aligned}$$

2. update velocity:

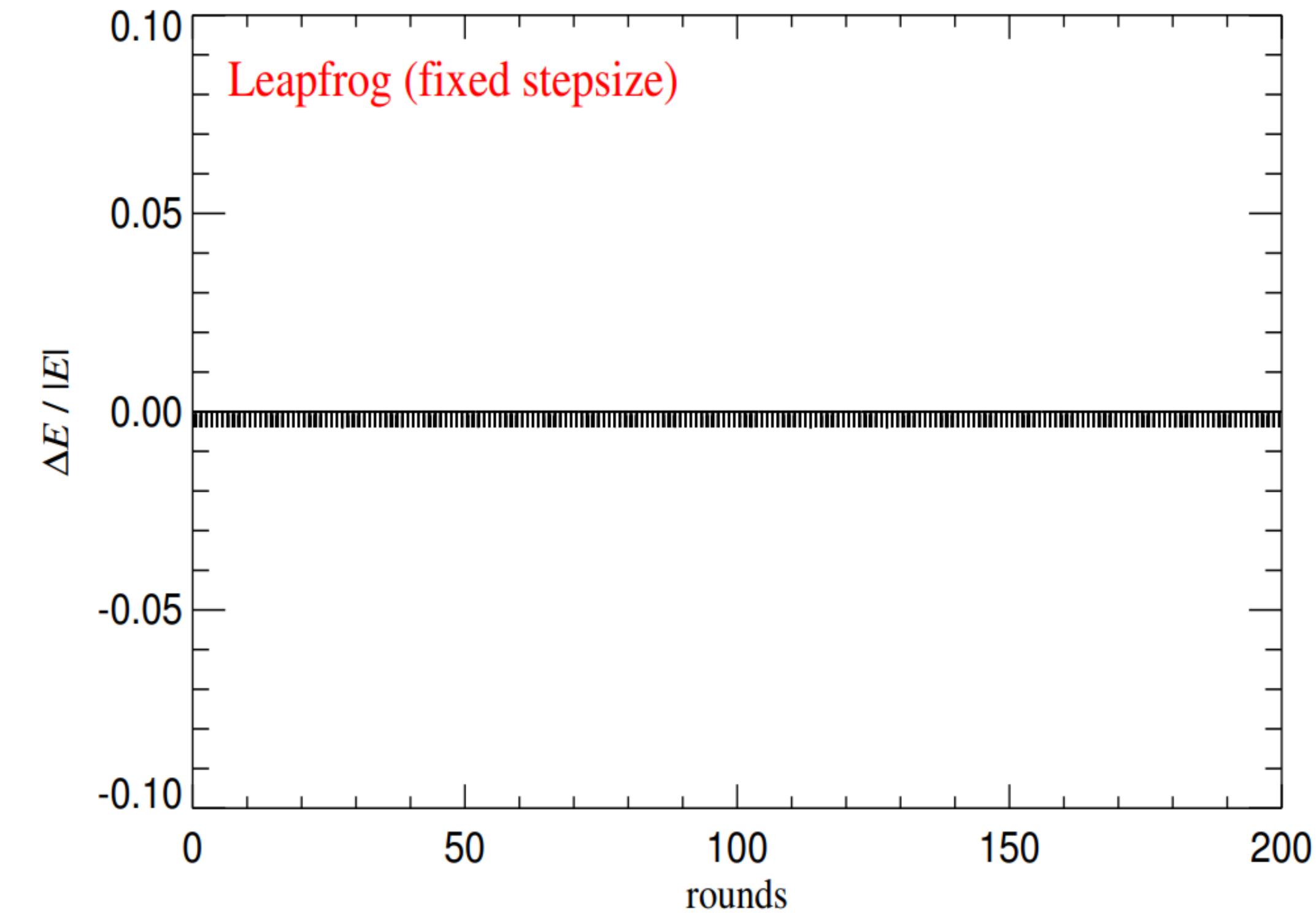
$$\begin{aligned}\underline{v}(t + \Delta t) = \underline{v}(t) + \\ \underline{a}\left(t + \frac{1}{2}\Delta t\right)\Delta t + \\ \mathcal{O}(\Delta t^3)\end{aligned}$$

SOLUTION OF THE TWO-BODY PROBLEM



Exact solution to $H_{\text{leap}} = H + H_{\text{err}}$,

$$H_{\text{err}} \propto \frac{\Delta t^2}{12} \left\{ \{H_{\text{kin}}, H_{\text{pot}}\}, H_{\text{kin}} + \frac{1}{2} H_{\text{pot}} \right\} + \mathcal{O}(\Delta t^3)$$



BREAK DOWN OF CLASSICAL METHODS

- Strong turbulence:
 - turbulent flow simulations are limited by Reynolds numbers. Unreasonable small grid required
 - → handle small scales by ML
- Problem of different scales
 - Different scales can be very problematic (although there are adaptive methods, mesh free methods, tree algorithms etc.)
- curse of dimensionality for high-dimensional PDEs
 - (e. g. Schrödinger in Physics, Black-Scholes in finance): The computational cost for solving them goes up exponentially with the dimensionality. [Han et al, 2018]
- Need for inverse modelling
- how to blend DEs with vast data sets?
- ...

Neural ODEs or Universal DEs

SCIENTIFIC ML — LEARNING THE SOLUTION OF ODES AND PDES

- XDEs are good for:
 - population models
 - motion of the planets
 - structural integrity of a bridge
 - fluid dynamics
 - ...

ODEs are kind of easy — only derivatives with respect to one variable

PDEs are more complicated — derivatives with respect to many variables and differential equations are local while solutions exhibit non-local properties

- Traditional solution: discretisation (in time and space) and iterative solution

NEURAL ODES, OPERATOR LEARNING AND PINNS

Neural ODE:

$$\frac{df}{dt} = h_\theta(x_0, t, p)$$

(neural net = right hand side of diff eq.
solution: integrate entire neural net.)

Neural Operator:

$$G_\theta : X \rightarrow Y \quad u \mapsto G_\theta(u) \text{ with } X, Y \text{ function spaces (infinite dimensional)}$$

(neural net approximates the operator
i.e. the map between function space)

PINN:

$$f(x, t) = h_\theta(x, t, p) \text{ with } \frac{df}{dt} = \frac{dh_\theta}{dt}, \quad \frac{df}{dx} = \frac{dh_\theta}{dx} \text{ need to fulfil the diff eq.}$$

(solution is given by neural net,
autodiff and diff eq. are used in loss)

An introduction to neural ODEs in scientific machine learning



Patrick Kidger



Cradle.bio

01

From traditional parameterised modelling to neural ODEs

ODEs are really really (really) good models.

$$\text{兔} (0) = 2 \quad \frac{d \text{兔}}{dt} = \text{兔} - \text{兔} \text{ 狐}$$
$$\text{狐} (0) = 1 \quad \frac{d \text{狐}}{dt} = -\text{狐} + \text{兔} \text{ 狐}$$

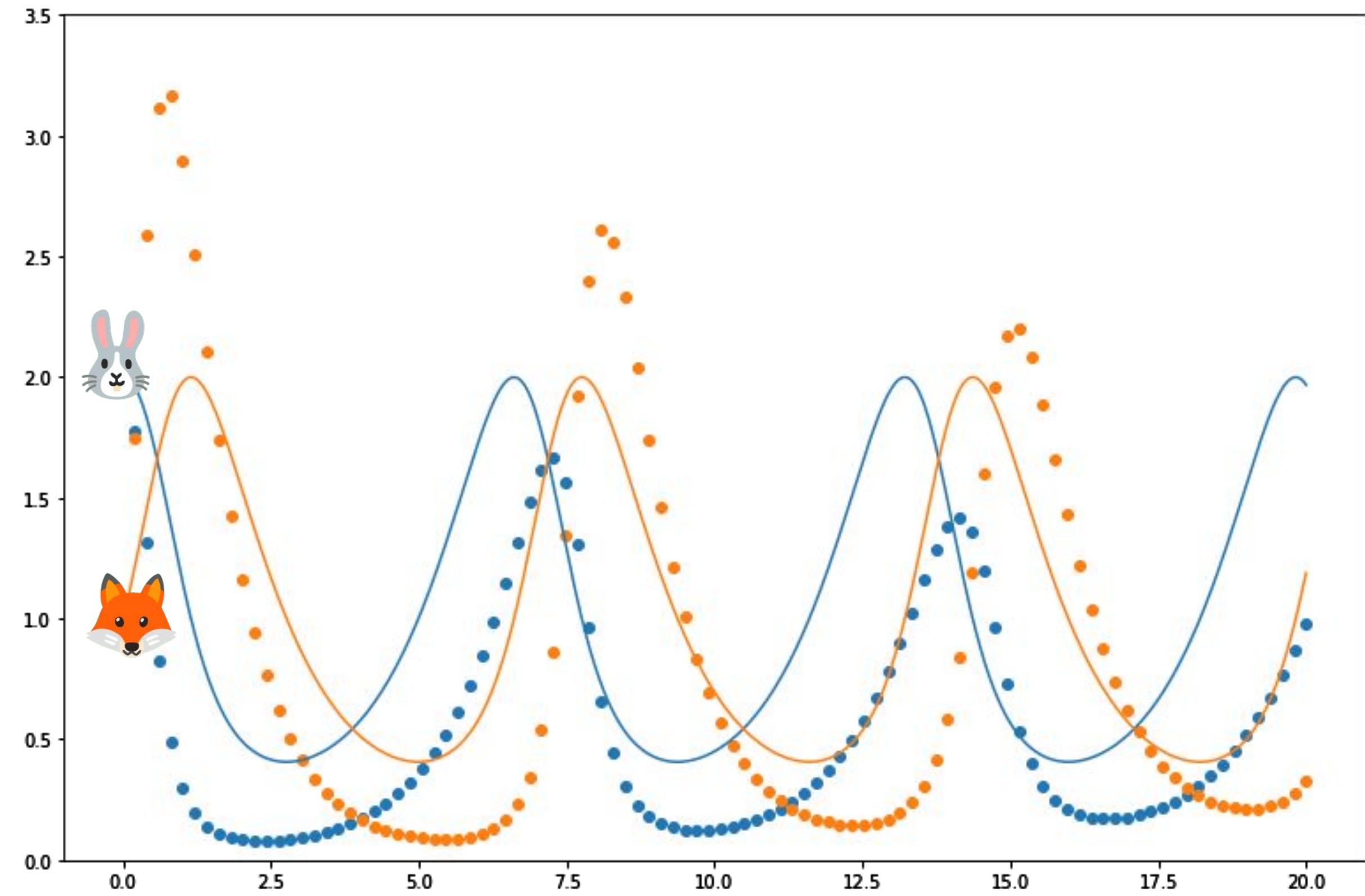
Good for:

- Population modelling;
- Motion of the planets;
- Structural integrity of a bridge;
- Fluid dynamics;
- ...

Simple idea: just write down how a value:

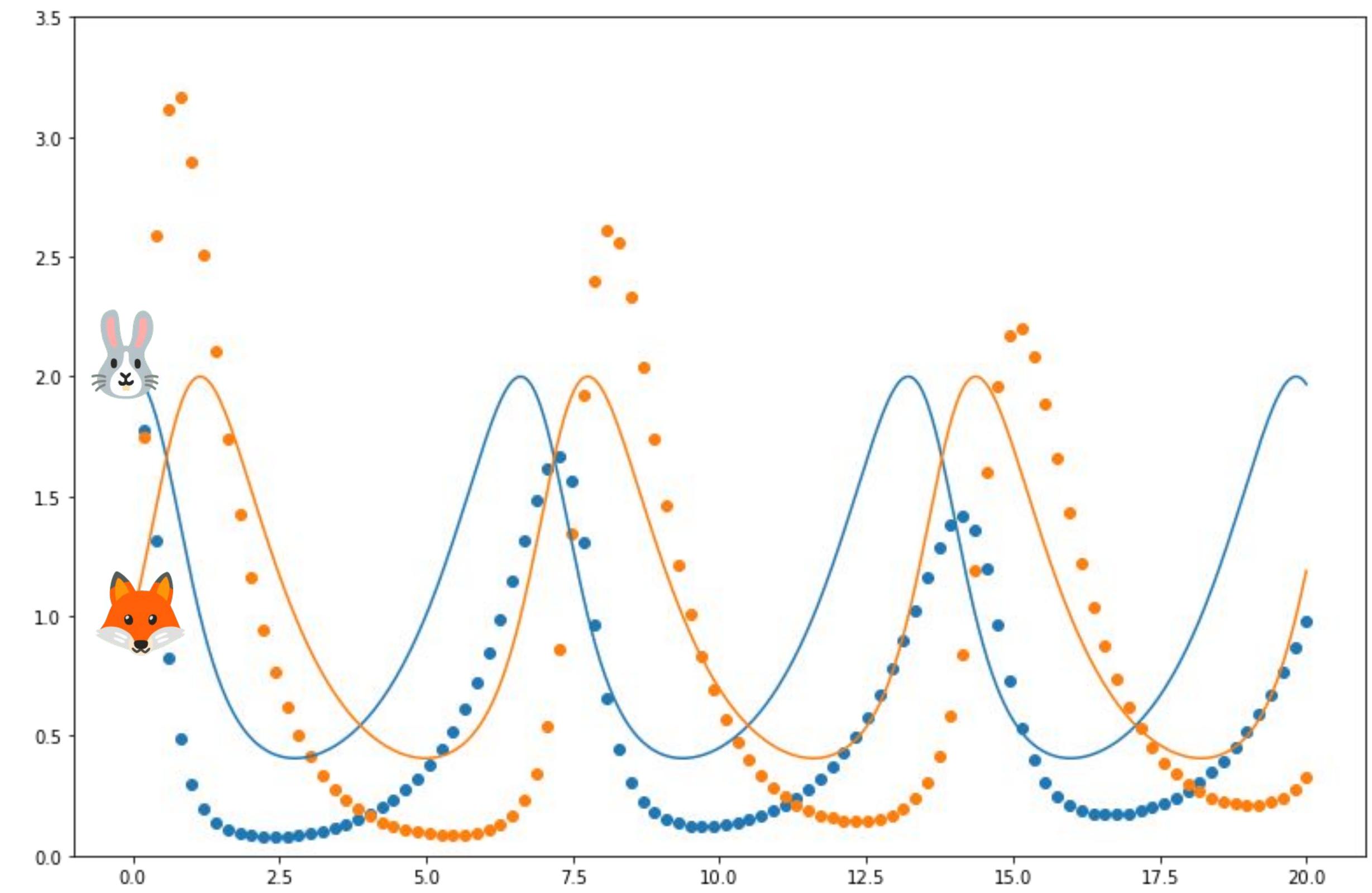
- Starts
- And how it changes over time.

Lotka–Volterra (Predator–Prey)



Parameterised ODEs

$$\text{🐰}(0) = 2 \quad \frac{d\text{🐰}}{dt} = \text{🐰} - \text{🐰} \text{🦊}$$
$$\text{🦊}(0) = 1 \quad \frac{d\text{🦊}}{dt} = -\text{🦊} + \text{🐰} \text{🦊}$$



Parameterised ODEs

- Initialise

$$a=1$$

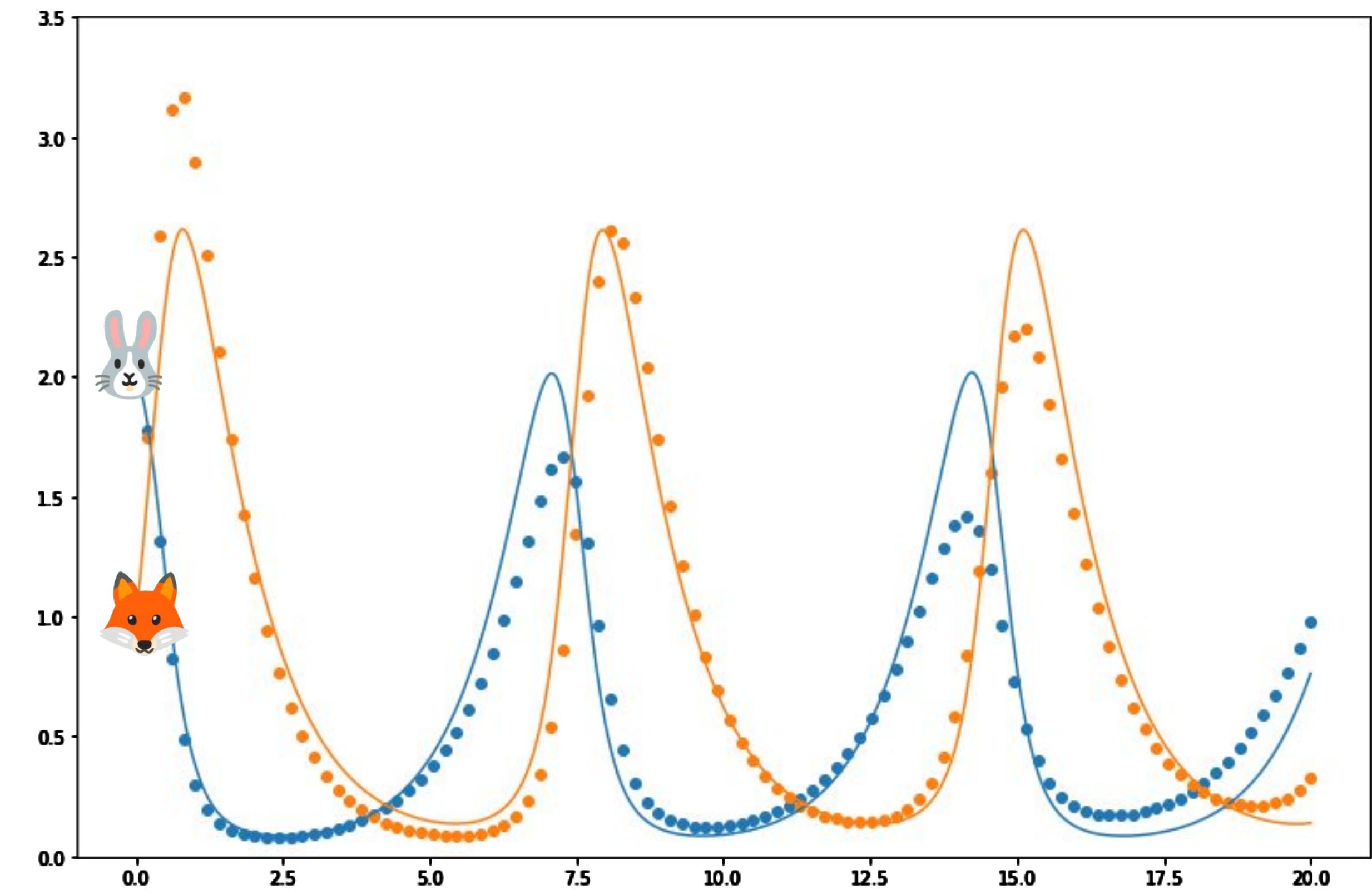
$$\beta=1$$

$$\gamma=1$$

$$\delta=1$$

- Set up a loss function:
 $(\text{data} - \text{model})^2$
- Backpropagate (through the ODE solver)
- Train via gradient descent

$$\begin{aligned} \text{🐰}(0) &= 2 & \frac{d\text{🐰}}{dt} &= a\text{🐰} - \beta\text{🐰}\text{🦊} \\ \text{🦊}(0) &= 1 & \frac{d\text{🦊}}{dt} &= -\gamma\text{🦊} + \delta\text{🐰}\text{🦊} \end{aligned}$$



Neural ODEs

- Initialise

$$\alpha=1$$

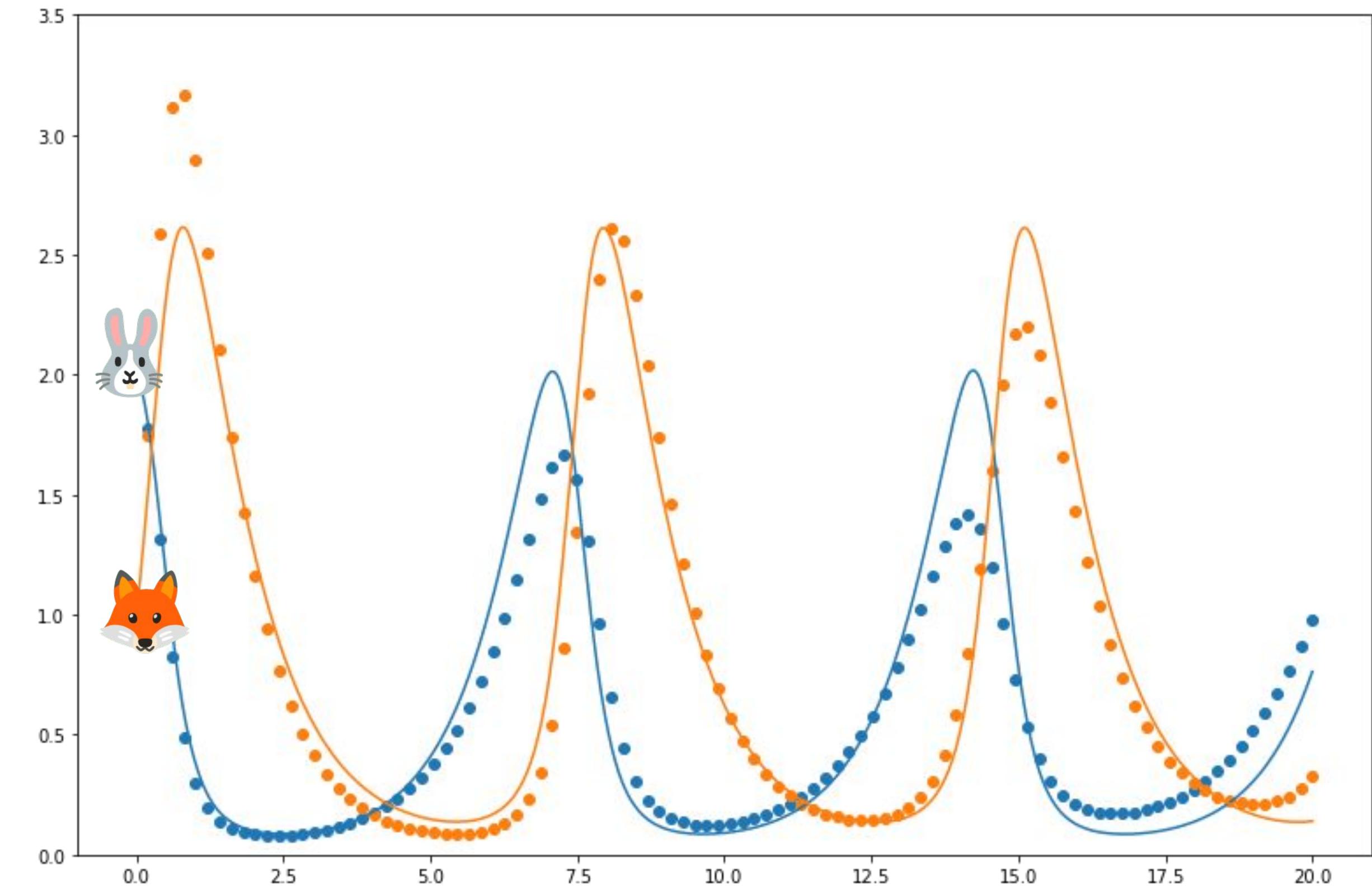
$$\beta=1$$

$$\gamma=1$$

$$\delta=1$$

- Set up a loss function:
 $(\text{data} - \text{model})^2$
- Backpropagate (through the ODE solver)
- Train via gradient descent

$$\begin{aligned} \text{兔}(0) &= 2 & \frac{d\text{兔}}{dt} &= \alpha \text{兔} - \beta \text{兔}\text{狐} + \text{NN}_{\theta}(\text{兔}, \text{狐}) \\ \text{狐}(0) &= 1 & \frac{d\text{狐}}{dt} &= -\gamma \text{狐} + \delta \text{兔}\text{狐} \end{aligned}$$



Neural ODEs

- Initialise

$$\alpha=1$$

$$\beta=1$$

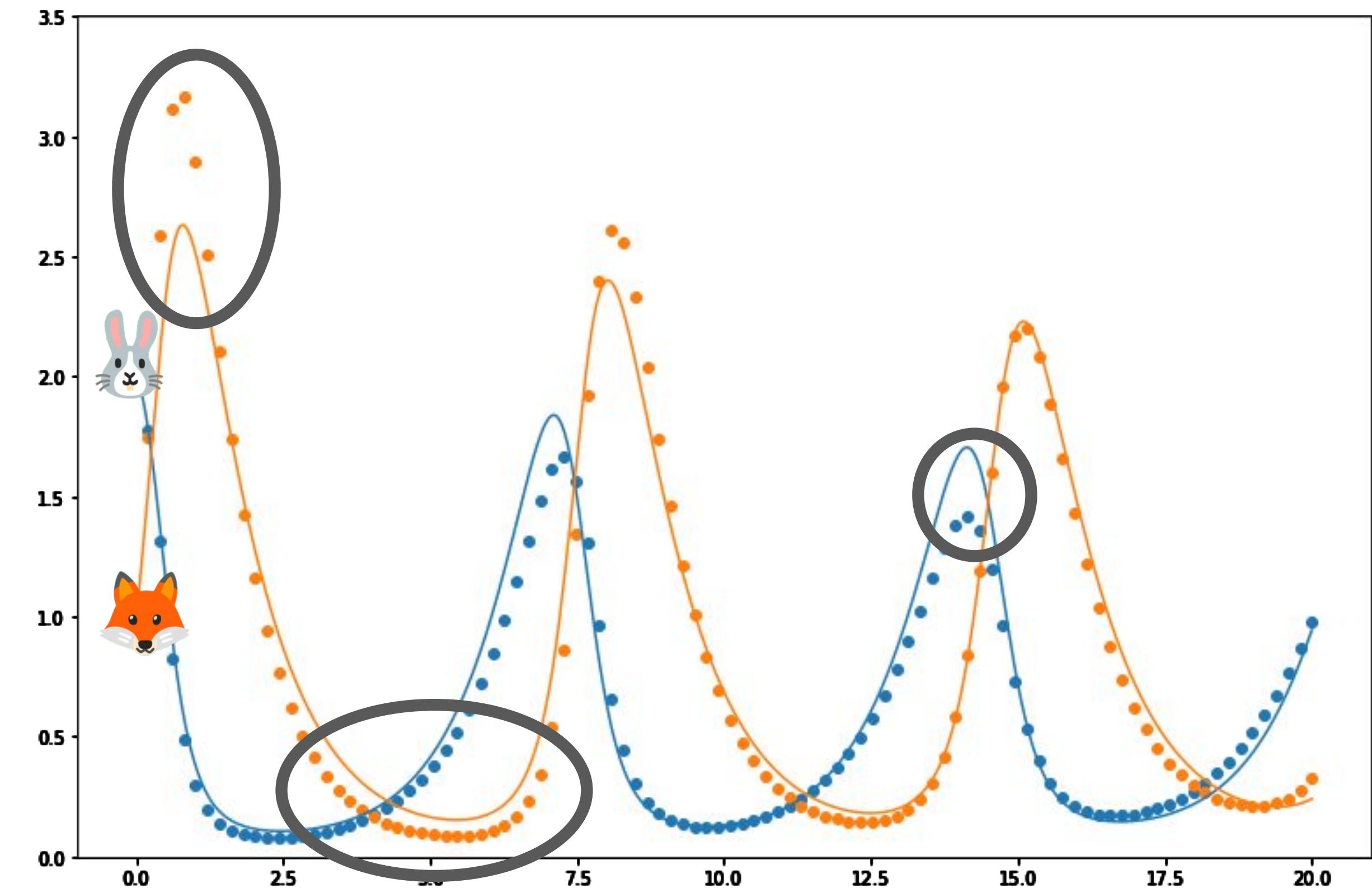
$$\gamma=1$$

$$\delta=1$$

$$\theta \sim N(0, \sigma^2)$$

- Set up a loss function:
 $(\text{data} - \text{model})^2$
- Backpropagate (through the ODE solver)
- Train via gradient descent

$$\begin{aligned} \text{🐰}(0) &= 2 & \frac{d\text{🐰}}{dt} &= \alpha\text{🐰} - \beta\text{🐰}\text{🦊} + \text{NN}_\theta(\text{🐰}, \text{🦊}) \\ \text{🦊}(0) &= 1 & \frac{d\text{🦊}}{dt} &= -\gamma\text{🦊} + \delta\text{🐰}\text{🦊} \end{aligned}$$



Neural ODEs: version 2

- Initialise

$$\alpha = 1$$

$$\beta = 1$$

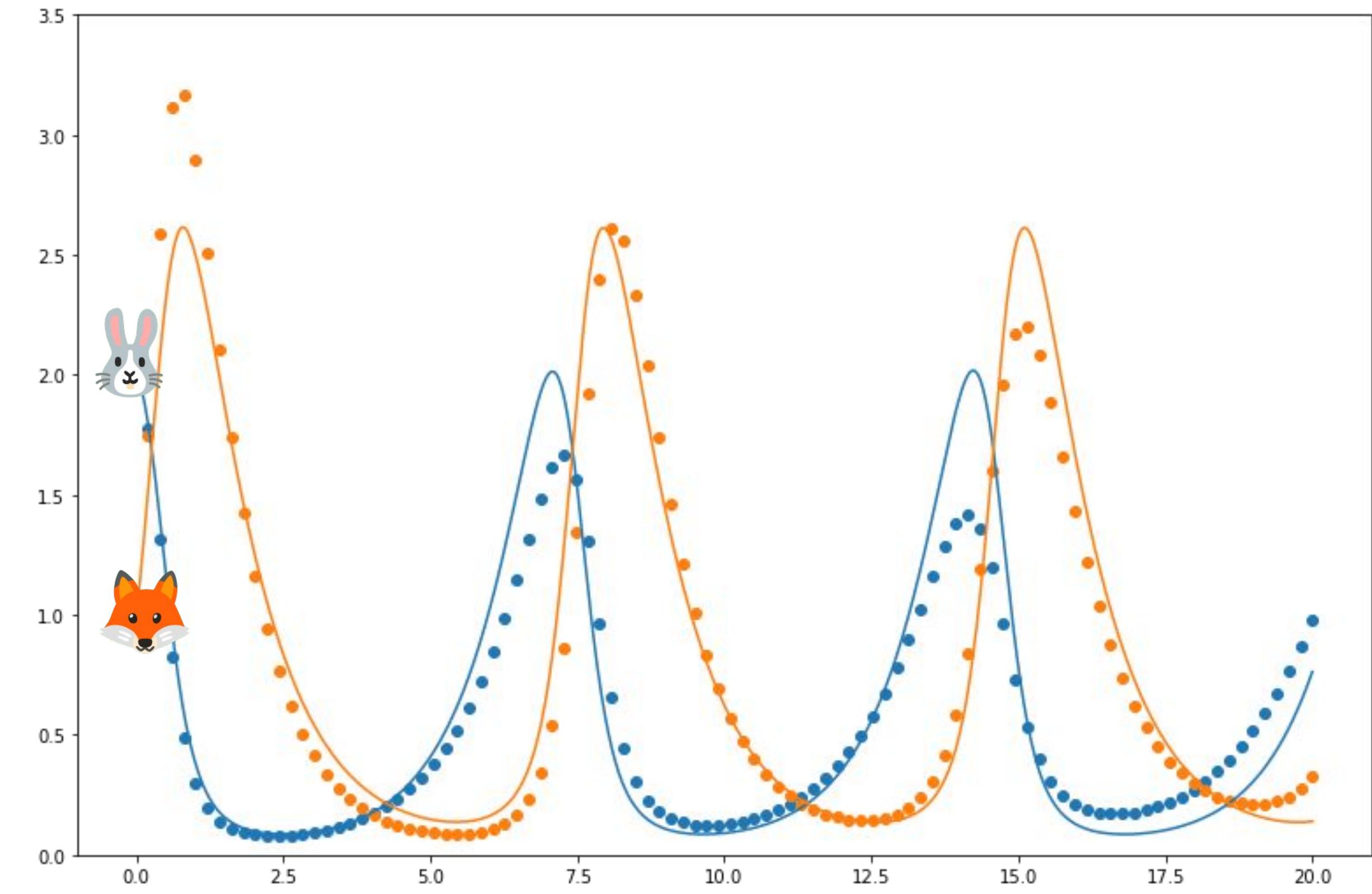
$$\gamma = 1$$

$$\delta = 1$$

$$\theta \sim N(0, \sigma^2)$$

- Set up a loss function:
 $(\text{data} - \text{model})^2$
- Backpropagate (through the ODE solver)
- Train via gradient descent

$$\begin{aligned} \text{🐰}(0) &= 2 & \frac{d\text{🐰}}{dt} &= \alpha \text{🐰} - \beta \text{🐰}\text{🦊} \\ \text{🦊}(0) &= 1 & \frac{d\text{🦊}}{dt} &= -\gamma \text{🦊} + \delta \text{🐰}\text{🦊} \end{aligned}$$



Neural ODEs: version 2

- Initialise

$$\alpha=1$$

$$\beta=1$$

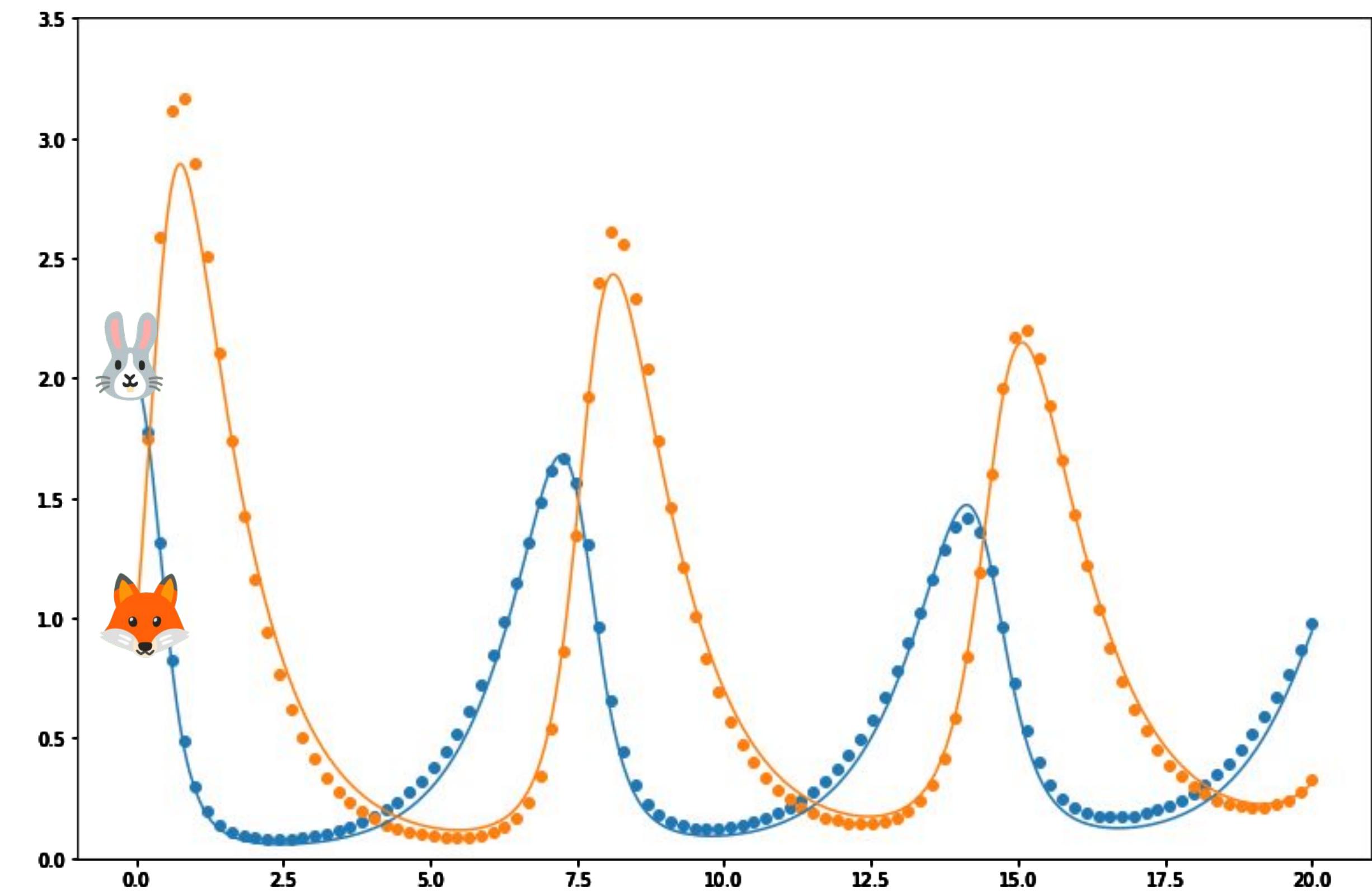
$$\gamma=1$$

$$\delta=1$$

$$\theta \sim N(0, \sigma^2)$$

- Set up a loss function:
 $(\text{data} - \text{model})^2$
- Backpropagate (through the ODE solver)
- Train via gradient descent

$$\begin{aligned} \text{🐰}(0) &= 2 & \frac{d\text{🐰}}{dt} &= \alpha\text{🐰} - \beta\text{🐰}\text{🦊} \\ \text{🦊}(0) &= 1 & \frac{d\text{🦊}}{dt} &= -\gamma\text{🦊} + \delta\text{🐰}(\text{🦊} + \text{NN}_\theta(\text{🐰}, \text{🦊})) \end{aligned}$$



Summary of part 1:

- Neural ODEs are a great model for unknown physics / unknown biology / ...
 - When you don't fully understand the physics.
 - When the physics is too expensive to fully simulate. (E.g. in weather.)
 - When a traditional parameterised model is difficult to optimise. (E.g. in biology.)
 - *Differentiable simulation, surrogate models.*
- Keep the physical model if you think it's good.
 - Put the NN on the bit you think isn't well-modelled.
 - Here the NN is tiny: 4 neurons wide, 1 layer deep.
 - Physical parameters trained first; them kept fixed whilst training NN.
- Neural ODE trained in the same way as the mechanistic ODE.
 - In both cases: a parameterised vector field trained with backprop+SGD.
 - Easy to implement: software for NN, ODEs, SGD, ...
- “Neural ODE” \approx “hybrid ODE” \approx “universal ODE” \approx “parameterised ODE”

Symbolic regression

Let's try to interpret NN_{θ} .

Generate a dataset of ~10k samples:

Symbolically regress y_i against $(\text{🐰}_i, \text{🦊}_i)$:

Retrain all parameters differentiably:

vs SINDy?

SINDy performs Lasso of $d(\text{🐰}, \text{🦊})/dt$ against $(\text{🐰}, \text{🦊})$. Thus it requires an estimate of the derivative.
That's not needed here!

Also, we're not constrained to just Lasso: other symbolic regression (e.g. PySR i.e. regularised evolution) can be used.

$$\begin{aligned} d\text{🐰}/dt &= 1.1\text{🐰} - 1.3\text{🐰}\text{🦊} \\ d\text{🦊}/dt &= -0.98\text{🦊} + 1.6\text{🐰} (\text{🦊} + NN_{\theta}(\text{🐰}, \text{🦊})) \end{aligned}$$

$\{(\text{🐰}_1, \text{🦊}_1, y_1), \dots, (\text{🐰}_n, \text{🦊}_n, y_n)\}$ with $y_i = NN_{\theta}(\text{🐰}_i, \text{🦊}_i)$.
(Sample 🐰_i and 🦊_i in your favourite way, e.g. points from typical ODE trajectories.)

{PySR, Lasso, ...} $\Rightarrow y = \text{🦊}^{0.9} - \text{🦊} + 0.5\text{🦊}^2$

$$\begin{aligned} d\text{🐰}/dt &= 0.9\text{🐰} - 1.1\text{🐰}\text{🦊} \\ d\text{🦊}/dt &= -1.2\text{🦊} + 2.1\text{🐰}\text{🦊}^{0.95} \end{aligned}$$

Conditions + Learnt Initialisations

Recall our general ODE formulation.

Both f and g may depend on additional data, e.g., the latitude L that the 🐰 and 🦊 live at.

Recall that an ODE is determined by:

- An initial condition;
- A vector field.

So far we've tried learning the vector field.

We can also learn the initial condition.

$$(\text{🐰, 🦊})(0) = g \quad \frac{d(\text{🐰, 🦊})}{dt} = f_{\theta}(\text{🐰, 🦊})$$

$$(\text{🐰, 🦊})(0) = g(L) \quad \frac{d(\text{🐰, 🦊})}{dt} = f_{\theta}(\text{🐰, 🦊}, L)$$

$$(\text{🐰, 🦊})(0) = g_{\theta}(L) \quad \frac{d(\text{🐰, 🦊})}{dt} = f_{\theta}(\text{🐰, 🦊}, L)$$

Generalised observables

We might observe things that aren't part of the evolving state of our ODE, but which still somehow give us information.

For example, we might observe the amount of rabbit... output 💩.

Now fit this to the data as well.

Even if we don't care about 💩, this can help us fit the 🐰, 🦊 parts better.

Important extreme case: "hidden state".

$$(\text{🐰}, \text{🦊})(0) = g_{\theta}(L) \quad \frac{d(\text{🐰}, \text{🦊})}{dt} = f_{\theta}(\text{🐰}, \text{🦊}, L)$$

$$\text{💩} = h_{\theta}(\text{🐰}, L)$$

$$\begin{aligned} \text{loss} = & \|\text{🐰} - \text{🐰}_{\text{data}}\|^2 \\ & + \|\text{🦊} - \text{🦊}_{\text{data}}\|^2 \\ & + \|\text{💩} - \text{💩}_{\text{data}}\|^2 \end{aligned}$$

$$\mathbf{y}(0) = g_{\theta}(c) \quad \frac{d\mathbf{y}}{dt} = f_{\theta}(\mathbf{y}, c) \quad (\text{🐰}, \text{🦊}, \text{💩}) = h_{\theta}(\mathbf{y}, c)$$

Mathematical Concepts of Neural ODEs

NEURAL ODES

- Traditional solution: discretisation (in time and space) and iterative solution

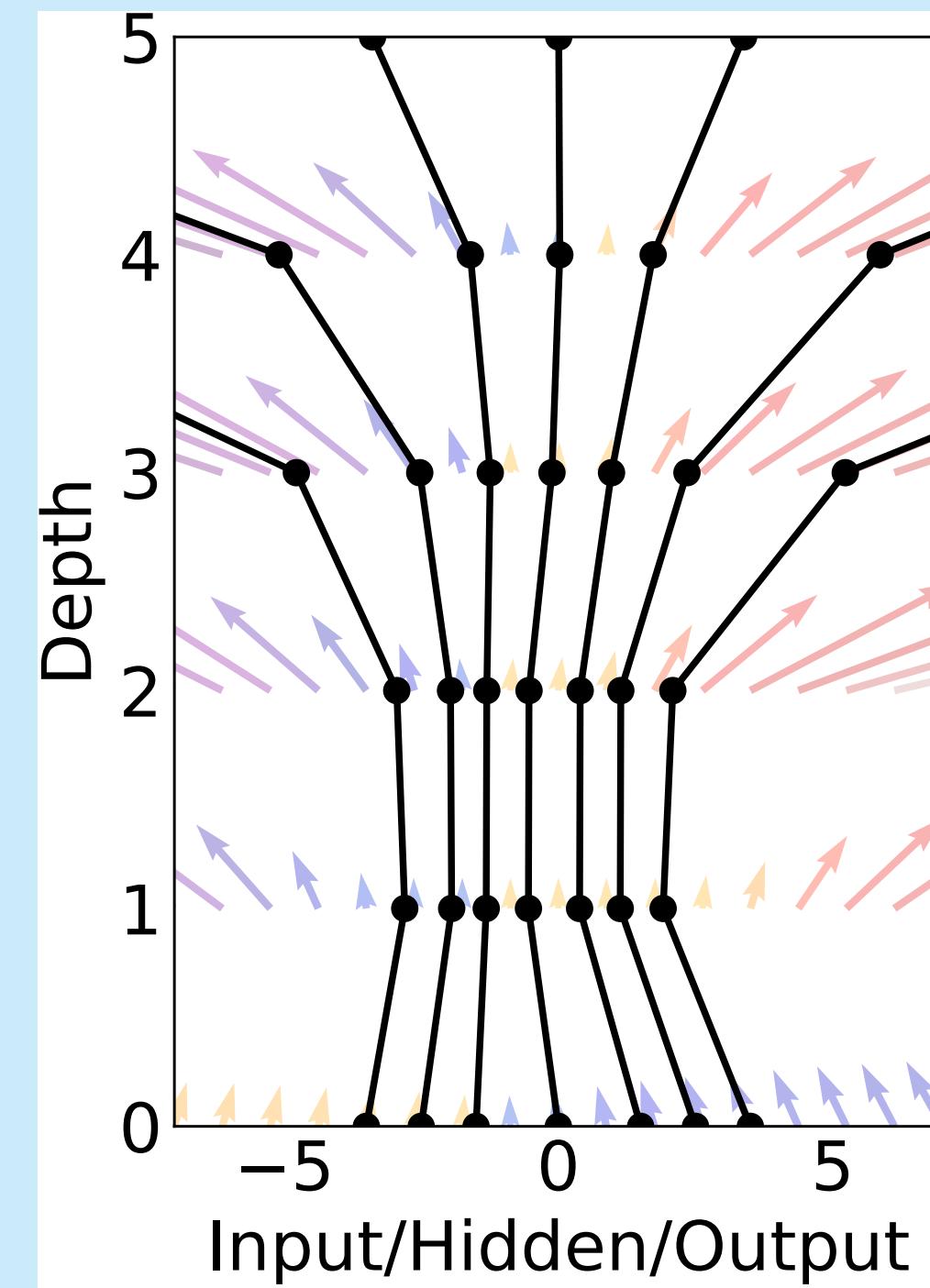
Euler discretization

$$h_{t+1} = h_t + f_t(h_t, \theta_t)$$

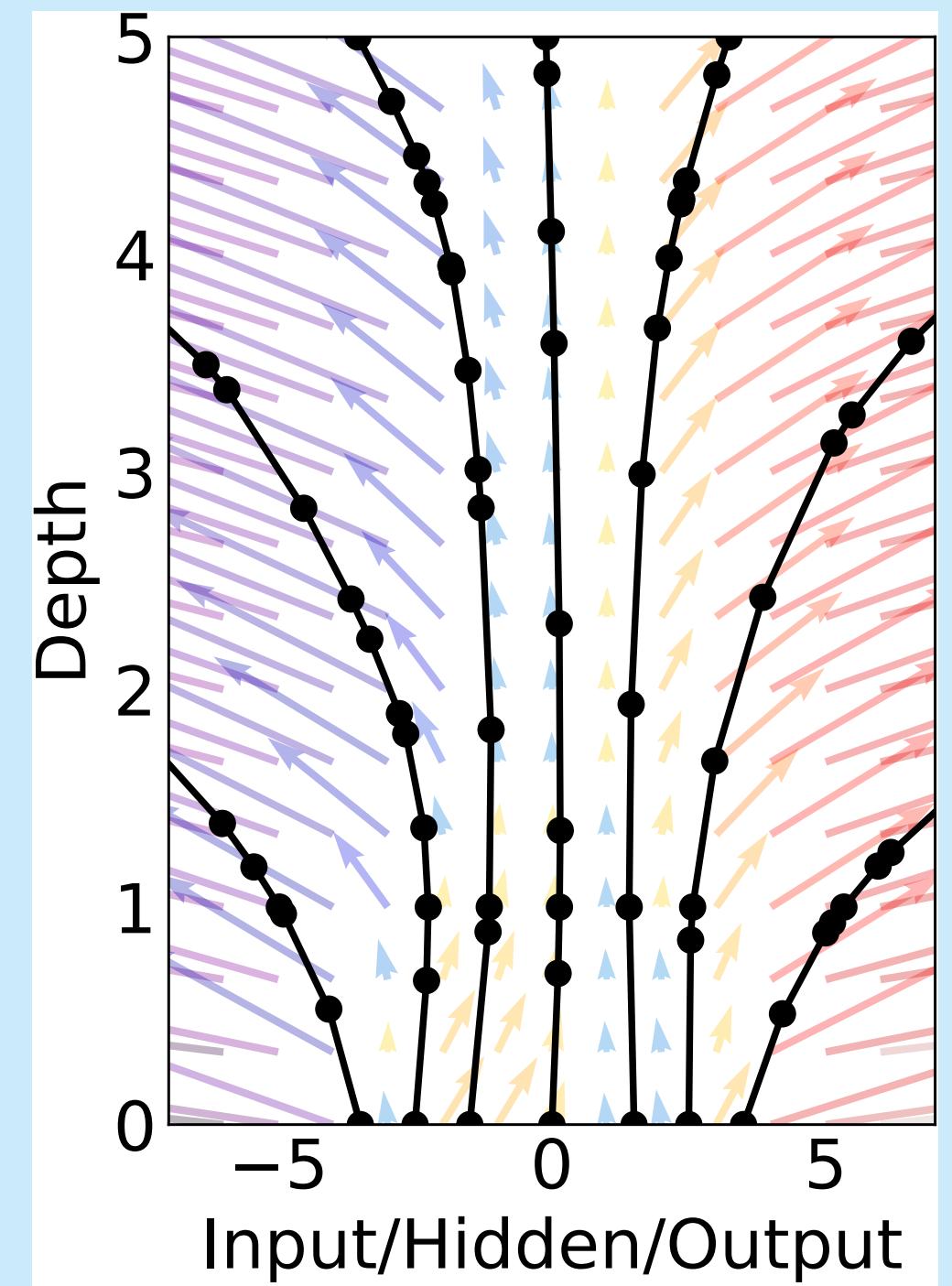
$$h(t + \Delta t) = h(t) + \Delta t \cdot f(t, h(t), \theta)$$

$$\frac{h(t + \Delta t) - h(t)}{\Delta t} = f(t, h(t), \theta)$$

Residual Network



ODE Network



Chen+2019

NEURAL ODES

- Traditional solution: discretisation (in time and space) and iterative solution

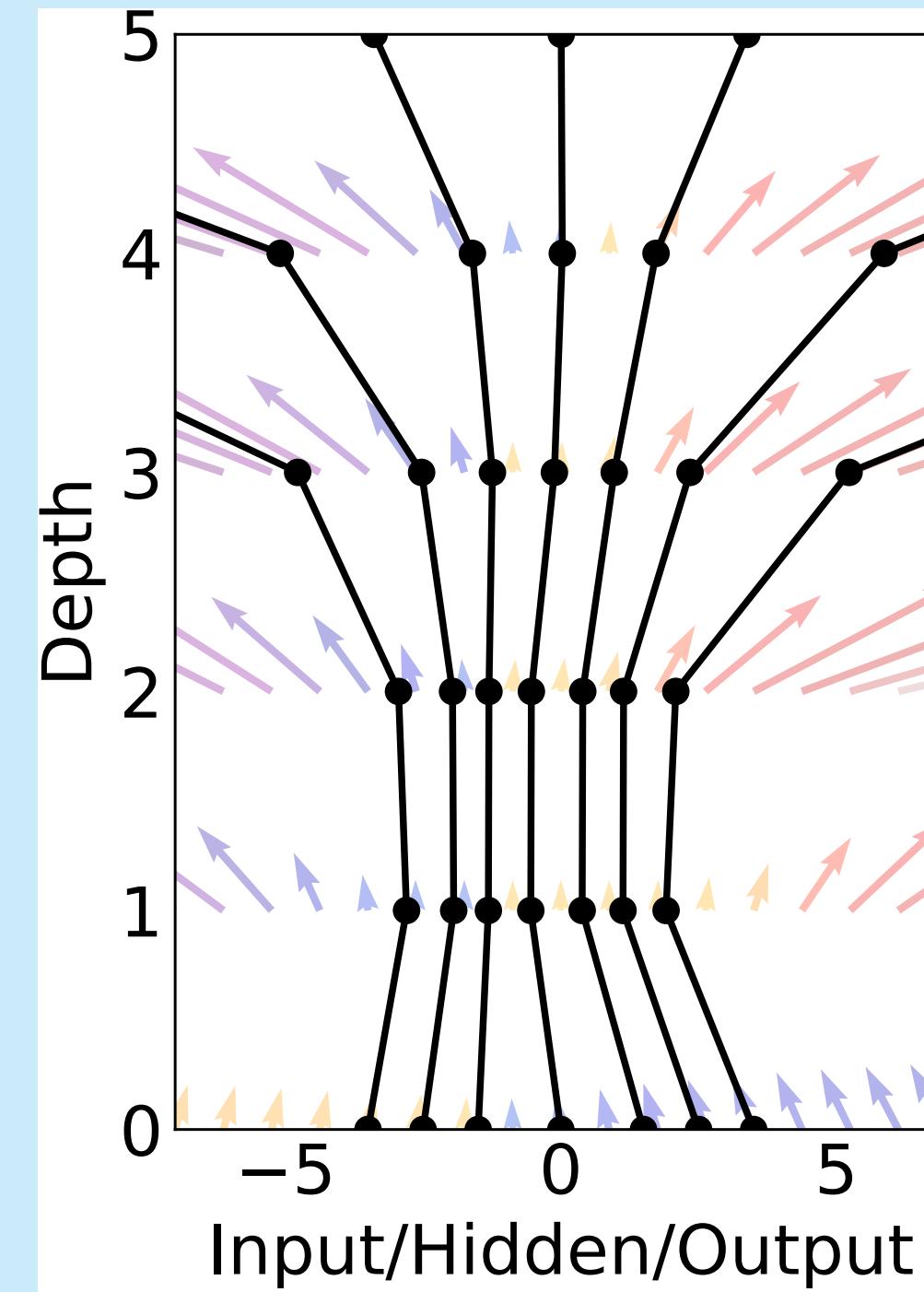
Euler discretization

$$h_{t+1} = h_t + f_t(h_t, \theta_t)$$

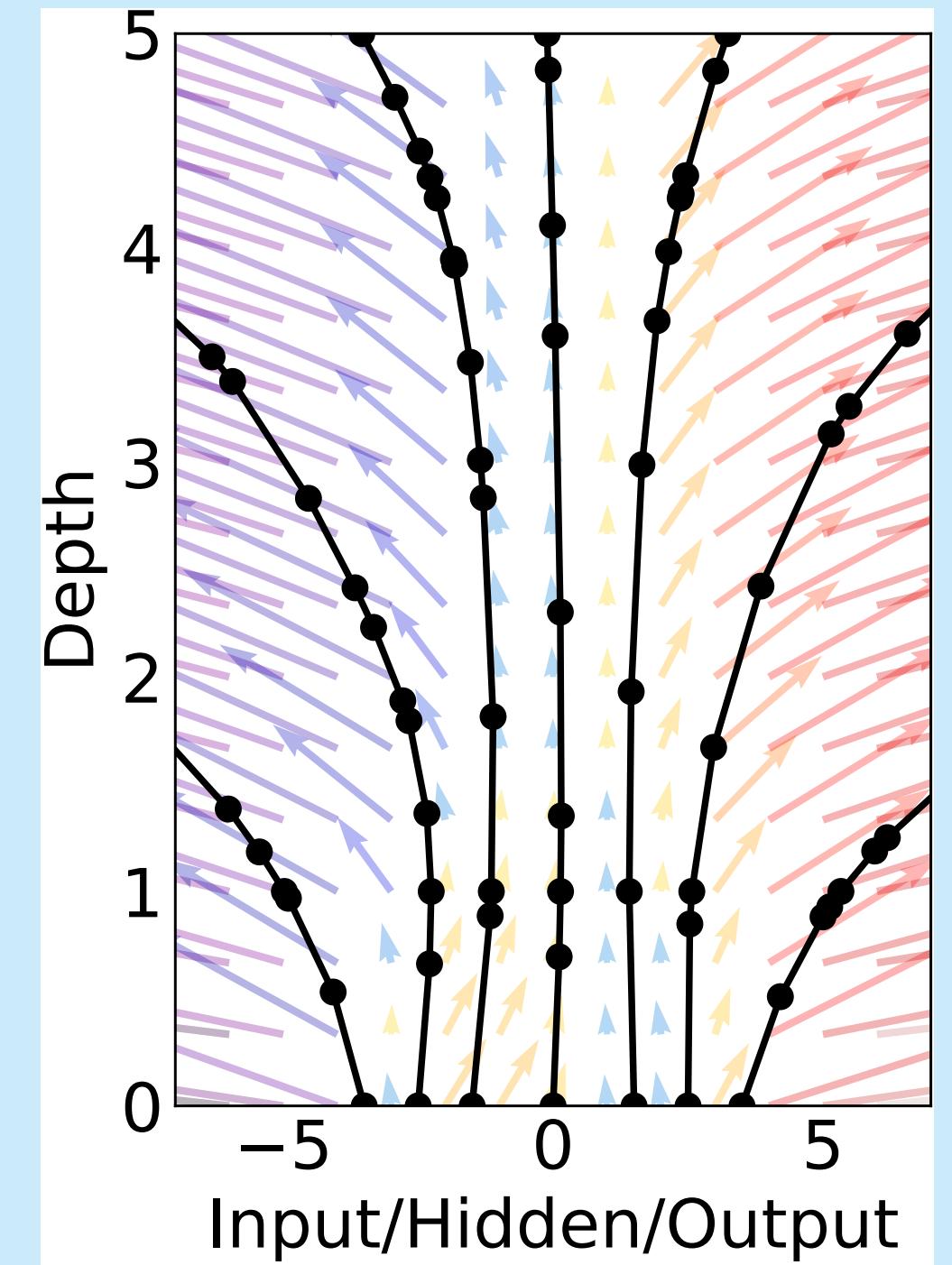
$$h(t + \Delta t) = h(t) + \Delta t \cdot f(t, h(t), \theta)$$

$$\frac{h(t + \Delta t) - h(t)}{\Delta t} = f(t, h(t), \theta)$$

Residual Network



ODE Network



Chen+2019

NEURAL ODES

What do we need to train a NeuralODE? → Reverse-mode automatic differentiation of ODE solutions

$$L(\mathbf{z}(t_1)) = L \left(\mathbf{z}(t_0) + \int_{t_0}^{t_1} f(\mathbf{z}(t), t, \theta) dt \right) = L(\text{ODESolve}(\mathbf{z}(t_0), f, t_0, t_1, \theta))$$

Problem is to find $\frac{dL(\mathbf{z}(t_1))}{d\theta}$ efficiently!

Don't worry about the details... diffrazx (Patrick Kidger), torchdiffeq (Ricky Chen), DiffEqFlux.jl got you convert!

NEURAL ODES TRAINING

In general two options for NODE training:

1. directly use autodiff for $\frac{dL(z(t_1))}{d\theta}$ and back propagate through the solver.
aka „discretize then optimise“

advantage: faster, more accurate gradient;

disadvantage: more memory consumption

2. solve $\frac{dL(z(t_1))}{d\theta}$ analytically.
Doing so will result in a backwards-in-time ODE (adjoint equation) that must be solved numerically
aka „optimise then discretise“

advantage: less memory consumption;

disadvantage: more computationally expensive, only
approximate gradients

ADJOINT EQUATION: PROBLEM STATEMENT

Find $\underset{\theta}{\operatorname{argmin}} \ L(z(t_1))$ (PM)

subject to

$$F(z(\dot{t}), z(t), \theta, t) = z(\dot{t}) - f(z(t), \theta, t) = 0 \quad (1)$$

$$z(t_0) = z_{t_0} \quad (2)$$

$$t_0 < t_1$$

- f is our neural network with parameters θ
- $z(t_0)$ is our input, $z(t_1)$ is the output, $z(t)$ is the state reached from $z(t_0)$ at time $t \in [t_0, t_1]$
- $\dot{z}(t) = \frac{dz(t)}{dt}$ is the time derivative of our state $z(t)$.
- L is our loss and its a function of the output $z(t_1)$.

ADJOINT EQUATION: PROBLEM STATEMENT

Find $\underset{\theta}{\operatorname{argmin}} \ L(z(t_1))$ (PM)

subject to

$$F(z(\dot{t}), z(t), \theta, t) = z(\dot{t}) - f(z(t), \theta, t) = 0 \quad (1)$$

$$z(t_0) = z_{t_0} \quad (2)$$

$$t_0 < t_1$$

- f is our neural network with parameters θ
- $z(t_0)$ is our input, $z(t_1)$ is the output, $z(t)$ is the state reached from $z(t_0)$ at time $t \in [t_0, t_1]$
- $\dot{z}(t) = \frac{dz(t)}{dt}$ is the time derivative of our state $z(t)$.
- L is our loss and its a function of the output $z(t_1)$.

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Let's rewrite our initial objective with a Lagrange multiplier $\lambda(t)$

$$\psi = L(z(t_1)) - \int_{t_0}^{t_1} \lambda(t) F(\dot{z}(t), z(t), \theta, t) dt$$

since

$$F(\dot{z}(t), z(t), \theta, t) = \dot{z}(t) - f(z(t), \theta, t) = 0$$

the integral vanishes and

$$\frac{d\psi}{d\theta} = \frac{dL(z(t_1))}{d\theta}$$

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Let's rewrite our initial objective with a Lagrange multiplier $\lambda(t)$

$$\psi = L(z(t_1)) - \int_{t_0}^{t_1} \lambda(t) F(\dot{z}(t), z(t), \theta, t) dt$$

since

$$F(\dot{z}(t), z(t), \theta, t) = \dot{z}(t) - f(z(t), \theta, t) = 0$$

the integral vanishes and

$$\frac{d\psi}{d\theta} = \frac{dL(z(t_1))}{d\theta}$$

Let's choose $\lambda(t)$ to eliminate hard to compute derivatives, such as the Jacobian $\frac{dz(t_1)}{d\theta}$

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Let the math magic happen...

(using integration by parts, chain rule and some recursive plugging in of equations...)
if you want the details, see the blog post.

$$\frac{dL}{d\theta} = \left[\frac{\partial L}{\partial z(t_1)} - \lambda(t_1) \right] \frac{dz(t_1)}{d\theta} + \int_{t_0}^{t_1} \left(\dot{\lambda}(t) + \lambda(t) \frac{\partial f}{\partial z} \right) \frac{dz(t)}{d\theta} dt + \int_{t_0}^{t_1} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Derivative Description	Ease
$\frac{\partial L}{\partial z(t_1)}$	Gradient of loss with respect to output. Easy
$\frac{dz(t_1)}{d\theta}$	Jacobian of output with respect to params. Not easy
$\dot{\lambda}(t)$	Derivative of lambda, a vector, with respect to time. Easy
$\lambda(t) \frac{\partial f}{\partial z}$	vector-Jacobian Product. Can compute with <u>reverse mode autodiff</u> without explicitly constructing Jacobian $\frac{\partial f}{\partial z}$. Easy
$\frac{dz(t)}{d\theta}$	Jacobian of arbitrary layer with respect to params. Not easy
$\lambda(t) \frac{\partial f}{\partial \theta}$	vector-Jacobian Product. Can compute with <u>reverse mode autodiff</u> without explicitly constructing Jacobian $\frac{\partial f}{\partial \theta}$. Easy

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Derivative Description	Ease
$\frac{\partial L}{\partial z(t_1)}$	Gradient of loss with respect to output. Easy
$\frac{dz(t_1)}{d\theta}$	Jacobian of output with respect to params. Not easy
$\dot{\lambda}(t)$	Derivative of lambda, a vector, with respect to time. Easy
$\lambda(t) \frac{\partial f}{\partial z}$	vector-Jacobian Product. Can compute with <u>reverse mode autodiff</u> without explicitly constructing Jacobian $\frac{\partial f}{\partial z}$. Easy
$\frac{dz(t)}{d\theta}$	Jacobian of arbitrary layer with respect to params. Not easy
$\lambda(t) \frac{\partial f}{\partial \theta}$	vector-Jacobian Product. Can compute with <u>reverse mode autodiff</u> without explicitly constructing Jacobian $\frac{\partial f}{\partial \theta}$. Easy

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Let the math magic happen...

(using integration by parts, chain rule and some recursive plugging in of equations...)
if you want the details, see the blog post.

$$\frac{dL}{d\theta} = \left[\frac{\partial L}{\partial z(t_1)} - \lambda(t_1) \frac{dz(t_1)}{d\theta} \right] + \int_{t_0}^{t_1} \left(\dot{\lambda}(t) + \lambda(t) \frac{\partial f}{\partial z} \right) \frac{dz(t)}{d\theta} dt + \int_{t_0}^{t_1} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

Can we get rid of $\frac{dz(t)}{d\theta}$ by a smart choice of $\lambda(t)$?

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

Let the math magic happen...

(using integration by parts, chain rule and some recursive plugging in of equations...)
if you want the details, see the blog post.

$$\frac{dL}{d\theta} = \underbrace{\left[\frac{\partial L}{\partial z(t_1)} - \lambda(t_1) \right]}_{=0} \underbrace{\frac{dz(t_1)}{d\theta}}_{=0} + \int_{t_0}^{t_1} \left(\dot{\lambda}(t) + \lambda(t) \frac{\partial f}{\partial z} \right) \underbrace{\frac{dz(t)}{d\theta}}_{=0} dt + \int_{t_0}^{t_1} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

$$\frac{dL}{d\theta} = \underbrace{\left[\frac{\partial L}{\partial z(t_1)} - \lambda(t_1) \right]}_{=0} \frac{dz(t_1)}{d\theta} + \int_{t_0}^{t_1} \underbrace{\left(\dot{\lambda}(t) + \lambda(t) \frac{\partial f}{\partial z} \right)}_{=0} \frac{dz(t)}{d\theta} dt + \int_{t_0}^{t_1} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

Let's define the backward ODE as:

$$\dot{\lambda}(t) = -\lambda(t) \frac{\partial f}{\partial z} \quad \text{s.t.} \quad \lambda(t_1) = \frac{\partial L}{\partial z(t_1)}$$

giving

$$\lambda(t_0) = \lambda(t_1) - \int_{t_1}^{t_0} \lambda(t) \frac{\partial f}{\partial z} dt$$

$\lambda(t)$ is called
the adjoint.

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

$$\frac{dL}{d\theta} = \underbrace{\left[\frac{\partial L}{\partial z(t_1)} - \lambda(t_1) \right]}_{=0} \frac{dz(t_1)}{d\theta} + \int_{t_0}^{t_1} \underbrace{\left(\dot{\lambda}(t) + \lambda(t) \frac{\partial f}{\partial z} \right)}_{=0} \frac{dz(t)}{d\theta} dt + \int_{t_0}^{t_1} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

Which simplifies our objective a lot!

$$\frac{dL}{d\theta} = - \int_{t_1}^{t_0} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

ADJOINT EQUATION VIA LAGRANGE MULTIPLIER

$$\frac{dL}{d\theta} = \underbrace{\left[\frac{\partial L}{\partial z(t_1)} - \lambda(t_1) \right]}_{=0} \frac{dz(t_1)}{d\theta} + \int_{t_0}^{t_1} \underbrace{\left(\dot{\lambda}(t) + \lambda(t) \frac{\partial f}{\partial z} \right)}_{=0} \frac{dz(t)}{d\theta} dt + \int_{t_0}^{t_1} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

Which simplifies our objective a lot!

$$\frac{dL}{d\theta} = - \int_{t_1}^{t_0} \lambda(t) \frac{\partial f}{\partial \theta} dt$$

(note, we switch the integration limits)

With this equation we trade memory for computation! Memory consumption is $O(1)$ with respect to the number of layers.

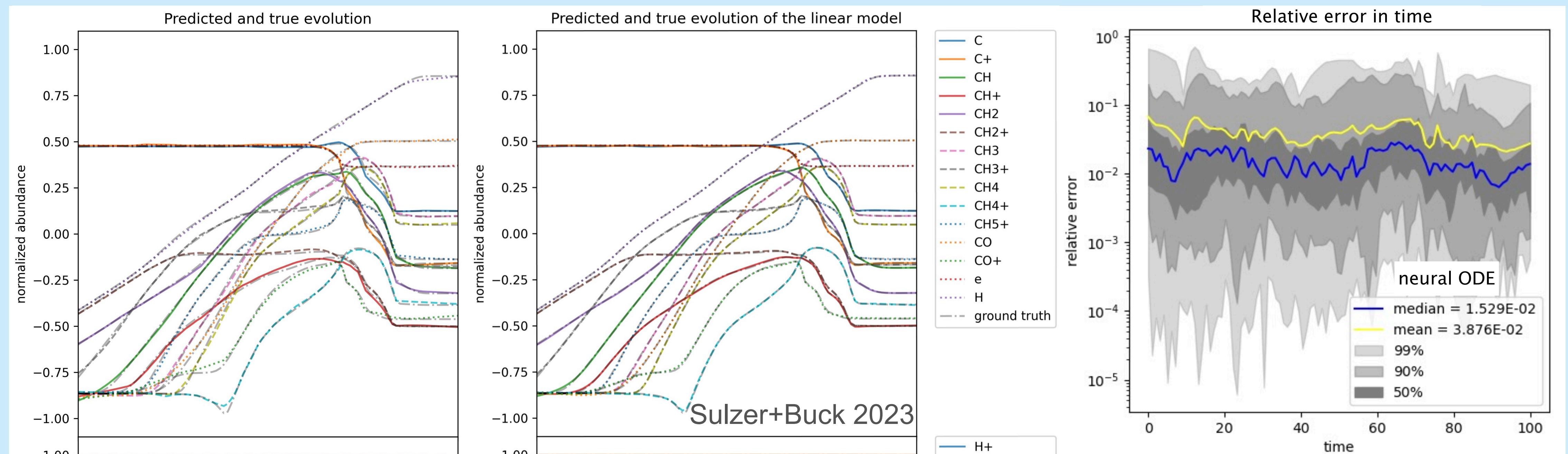
TRAINING NODE WITH ADJOINT METHOD

STEP BY STEP RECIPE

1. Forward pass: Solve the ODE (1)-(2) from time t_0 to t_1 and get the output $z(t_1)$.
2. Loss calculation: Calculate $L(z(t_1))$.
3. Backward pass: Solve ODEs (11) and (12) from reverse time t_1 to t_0 to get the gradient of the loss $\frac{dL(z(t_1))}{d\theta}$.
4. Use the gradient to update the network parameters θ .

NEURAL ODES IN ASTROPHYSICS

- Neural Astrophysical Wind Models (Nguyen 2023)
- Neural ODEs as a discovery tool to characterize the structure of the hot galactic wind of M82 (Nguyen+2023)
- Speeding up astrochemical reaction networks with autoencoders and neural ODEs (Sulzer+Buck 2023)



SUMMARY & CONCLUSION

Main take away:

scientific motivated inductive bias helps to be more robust, more data efficient
and better interpretable

My personal message:

Write better code!
Share more data!
Build more open-source software!

This will accelerate research cycles and lets you engage with peers early on!

WHERE TO FIND THE LECTURE MATERIAL



<https://github.com/TobiBu/graddays>

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.

