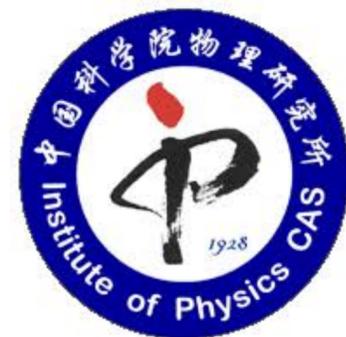


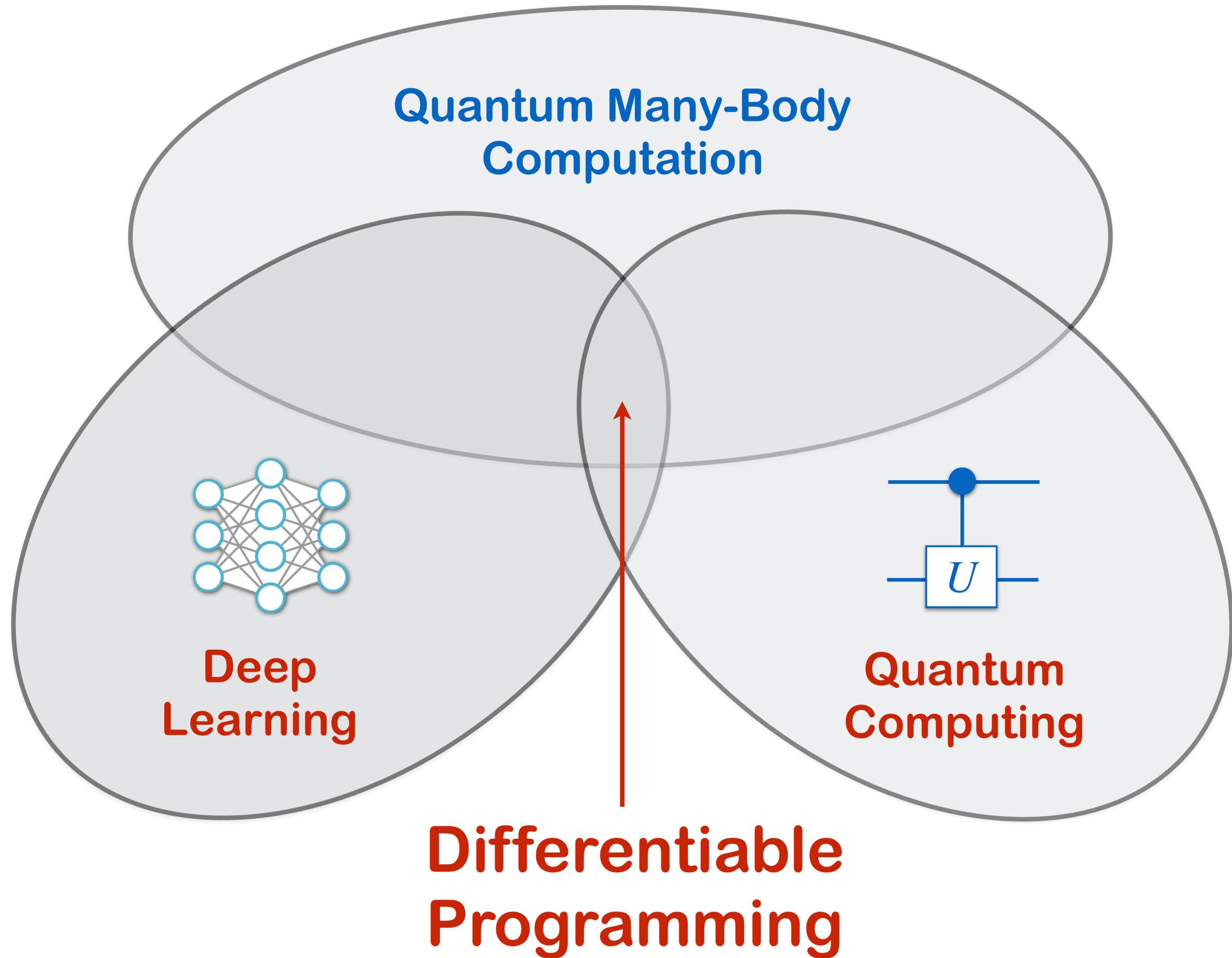
*d*ifferentiate everything: A lesson from deep learning

Lei Wang (王磊)

<https://wangleiphy.github.io>

Institute of Physics, CAS





Differentiable Programming

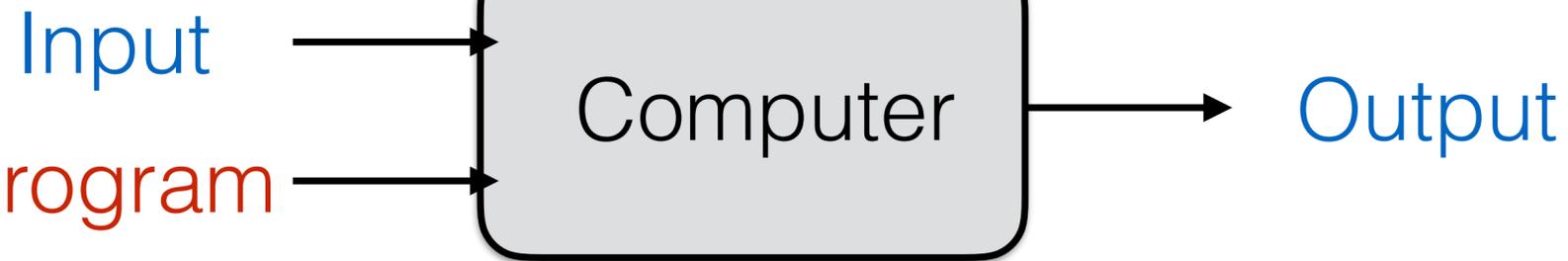


Andrej Karpathy

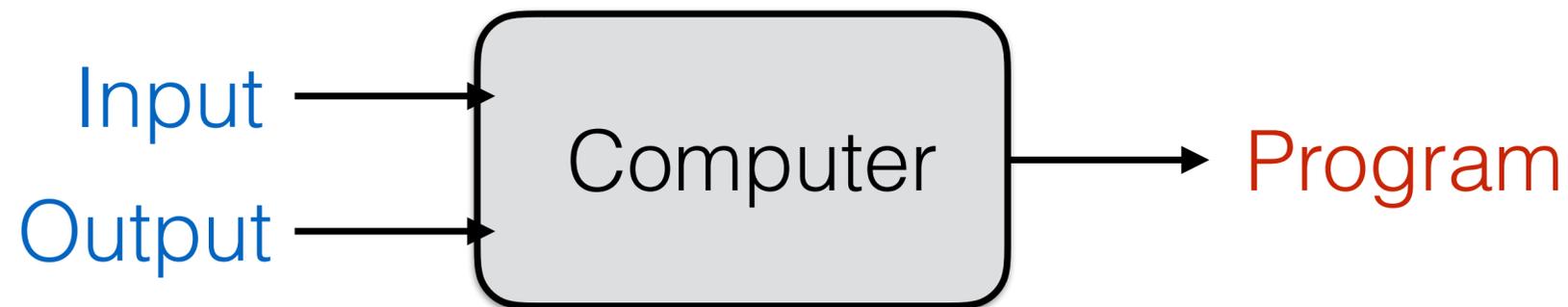
Director of AI at Tesla. Previously Research Scientist at OpenAI and PhD student at Stanford. I like to train deep neural nets on large datasets.

<https://medium.com/@karpathy/software-2-0-a64152b37c35>

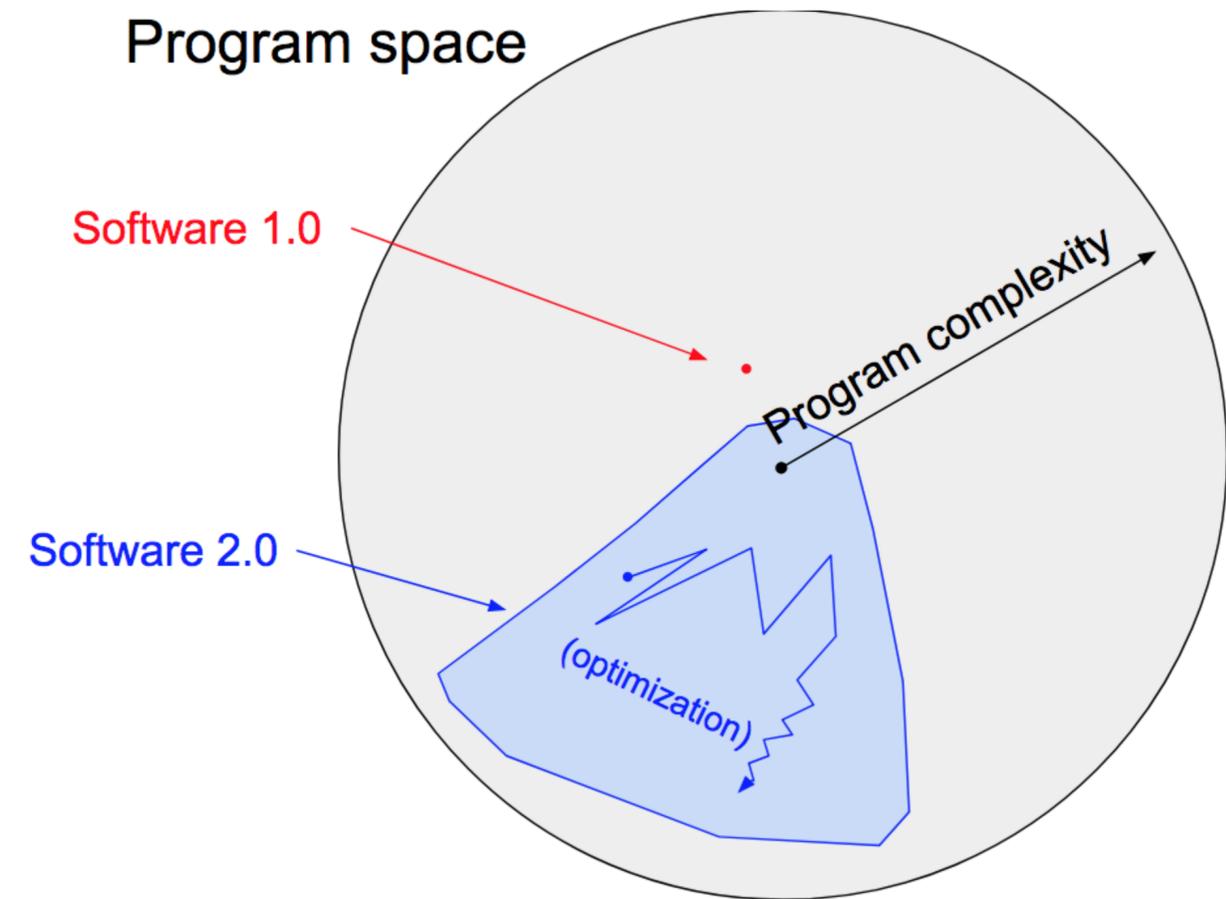
Traditional



Machine Learning



Program space



Writing software 2.0 by gradient search in the program space

Differentiable Programming

Benefits of Software 2.0

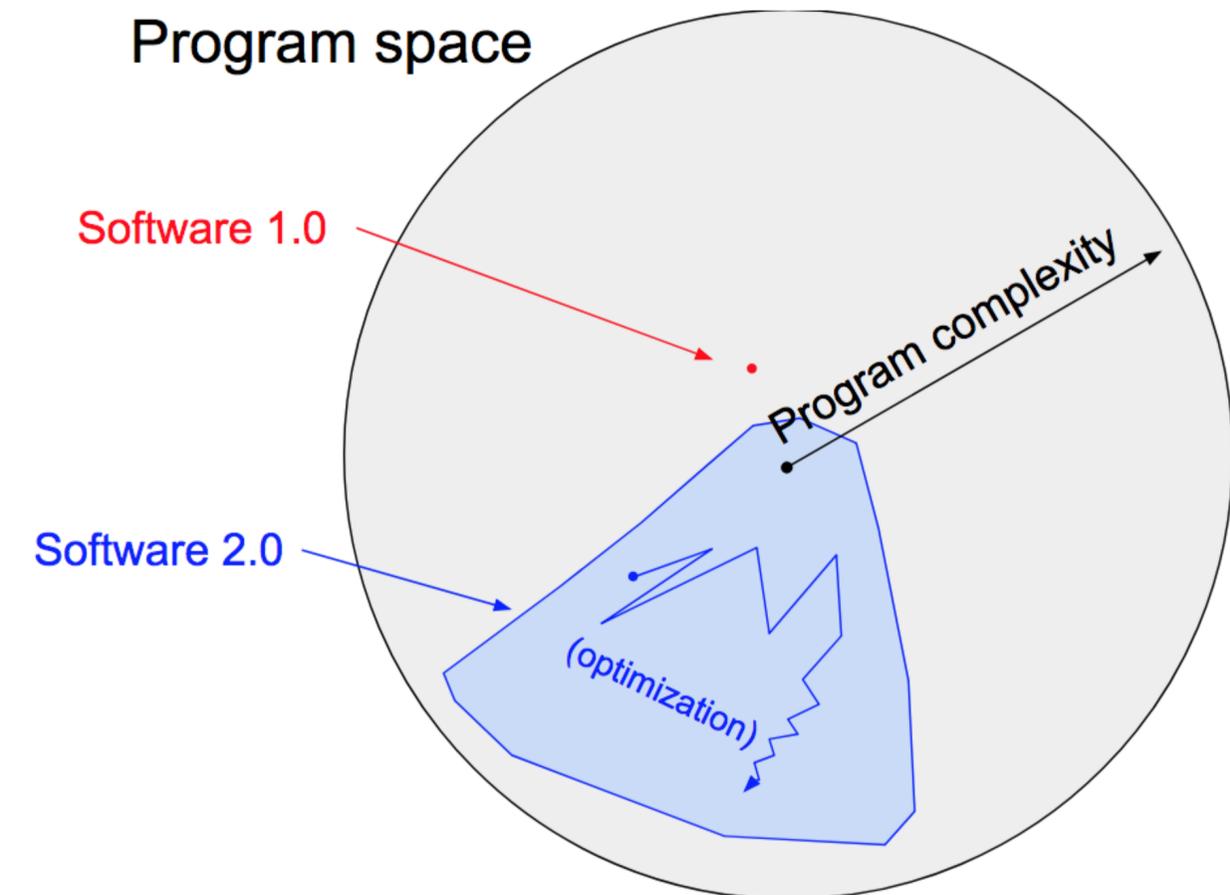
- Computationally homogeneous
- Simple to bake into silicon
- Constant running time
- Constant memory usage
- Highly portable & agile
- Modules can meld into an optimal whole
- **Better than humans**



Andrej Karpathy

Director of AI at Tesla. Previously Research Scientist at OpenAI and PhD student at Stanford. I like to train deep neural nets on large datasets.

<https://medium.com/@karpathy/software-2-0-a64152b37c35>

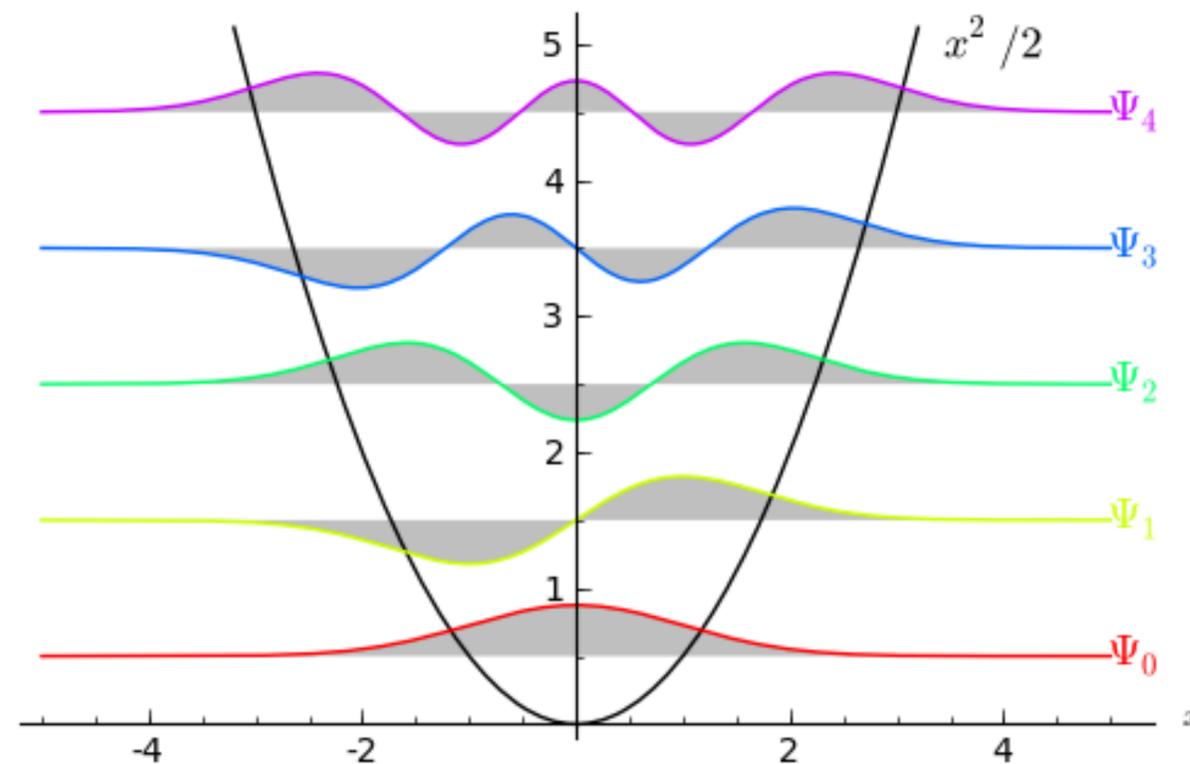


Writing software 2.0 by gradient search in the program space

Demo: Inverse Schrodinger Problem

Given ground state density, how to design the potential ?

$$\left[-\frac{1}{2} \frac{\partial^2}{\partial x^2} + V(x) \right] \Psi(x) = E \Psi(x)$$

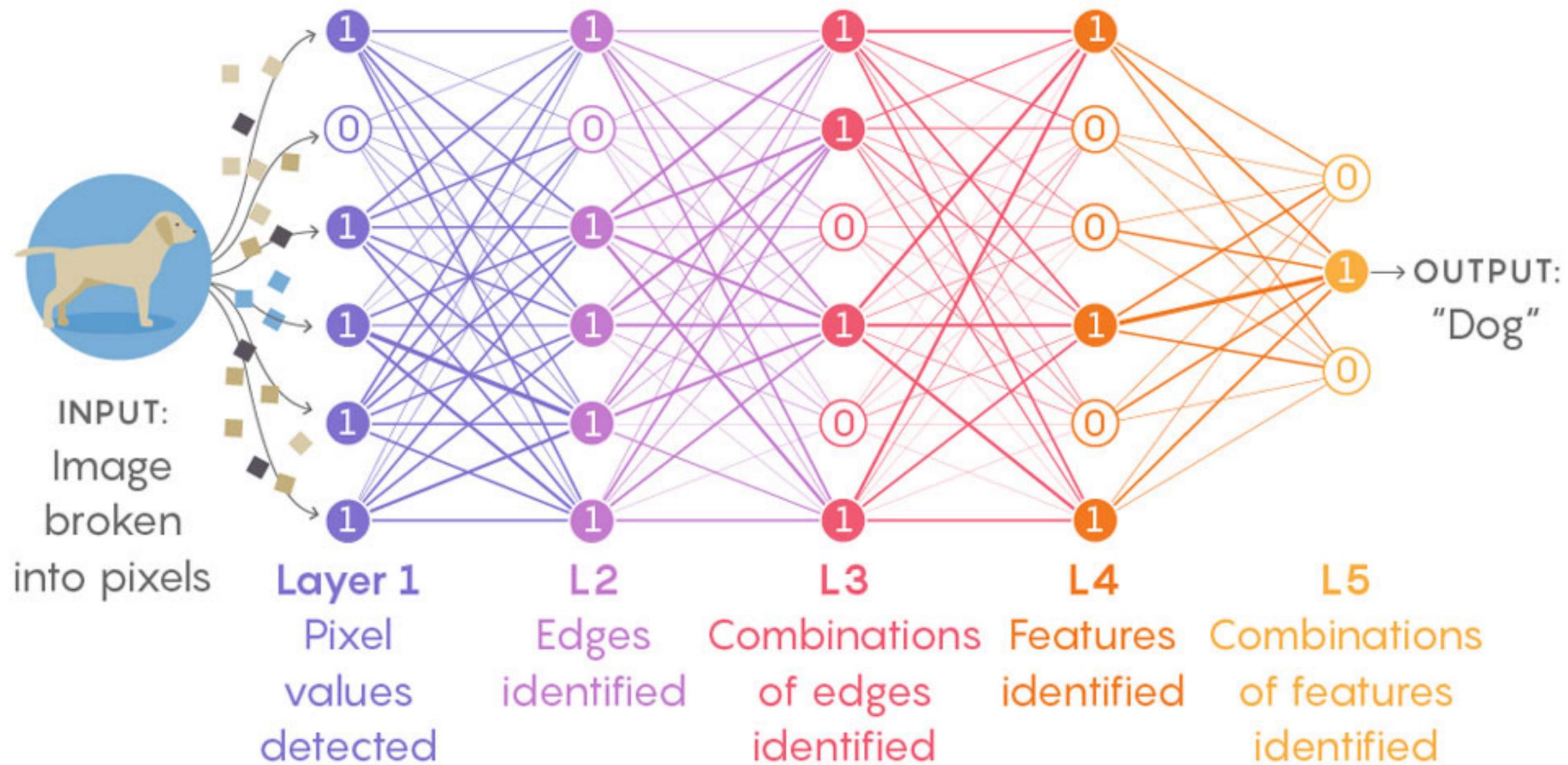


<https://math.mit.edu/~stevenj/18.336/adjoint.pdf>

https://github.com/QuantumBFS/SSSS/blob/master/1_deep_learning/schrodinger.py

What is under the hood?

What is deep learning ?

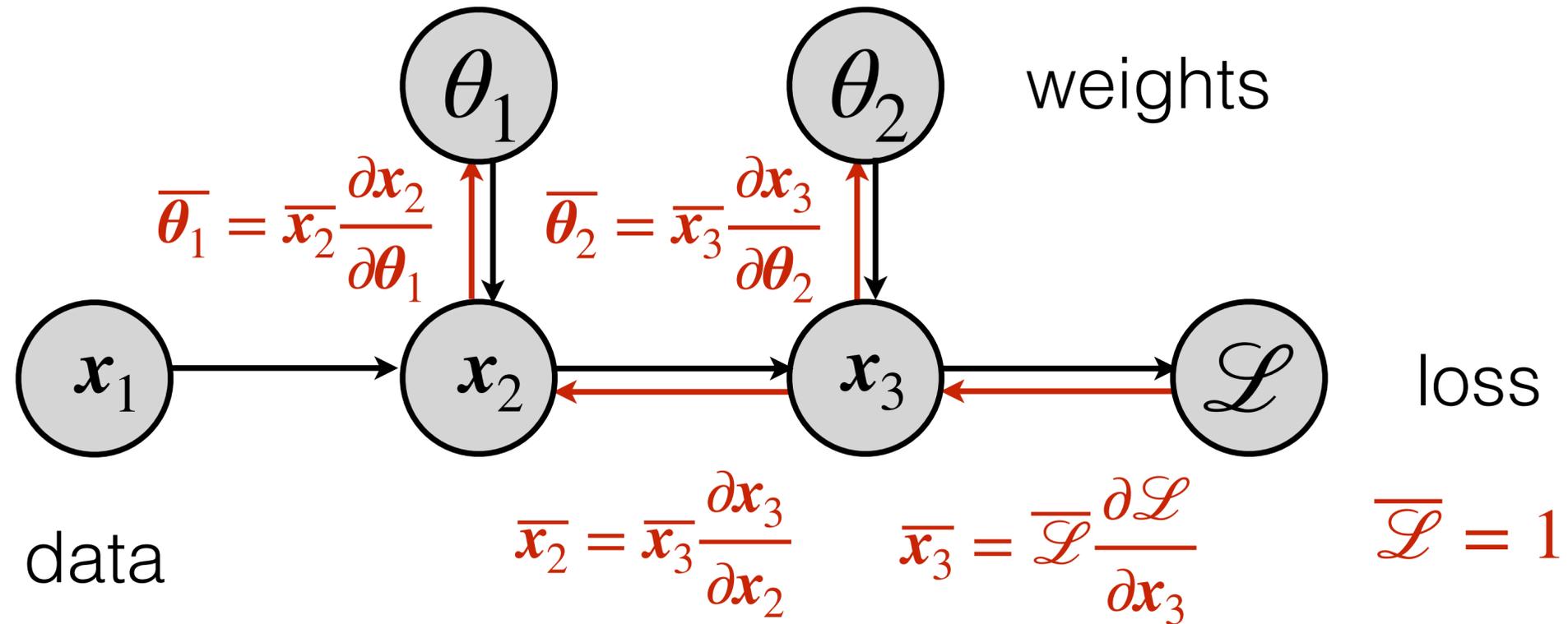


Composes differentiable components to a program e.g. a neural network, then optimizes it with gradients

Automatic differentiation on computation graph



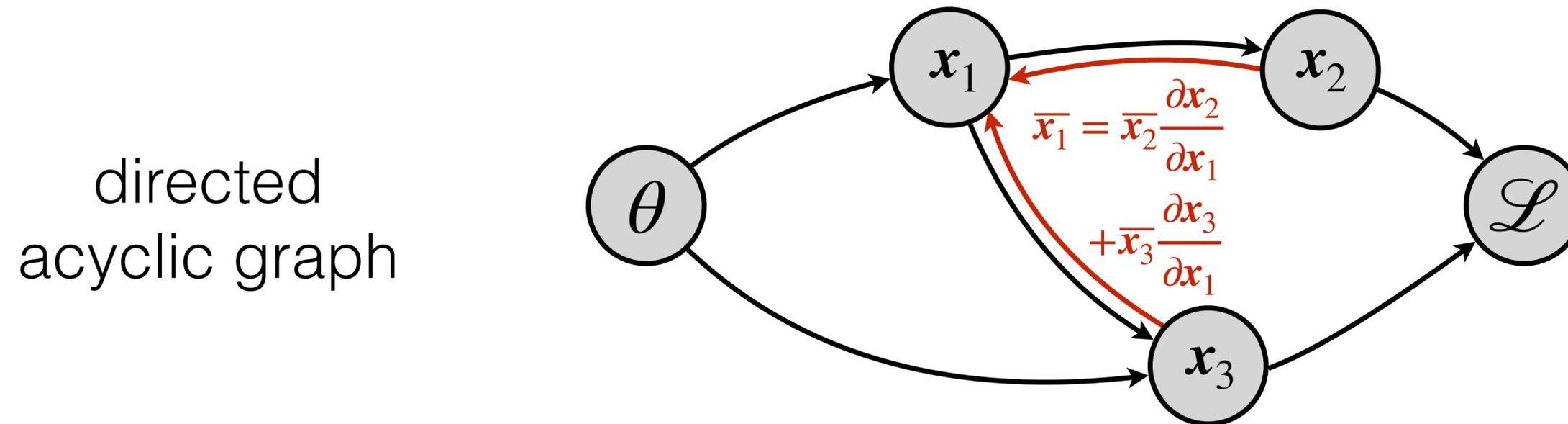
“comb graph”



“adjoint variable” $\bar{x} = \frac{\partial \mathcal{L}}{\partial x}$

Pullback the adjoint through the graph

Automatic differentiation on computation graph

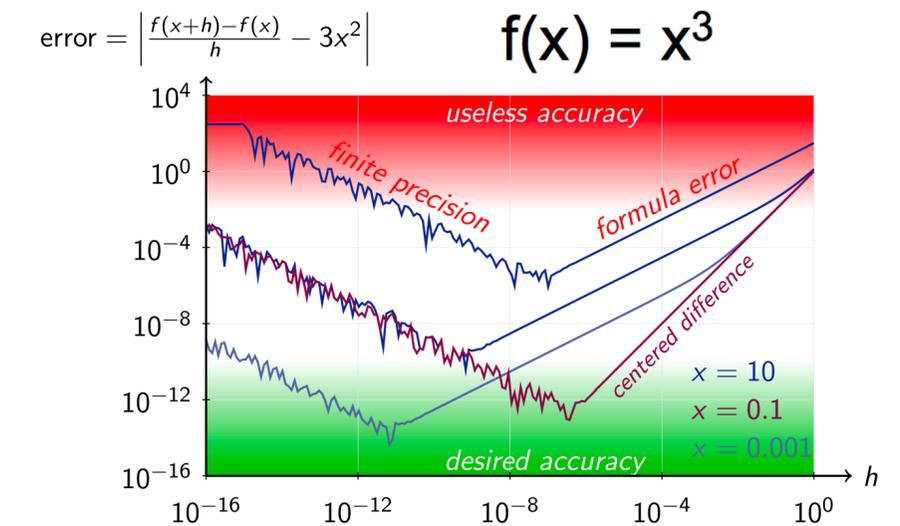


$$\bar{x}_i = \sum_{j: \text{child of } i} \bar{x}_j \frac{\partial x_j}{\partial x_i} \quad \text{with } \bar{L} = 1$$

Message passing for the adjoint at each node

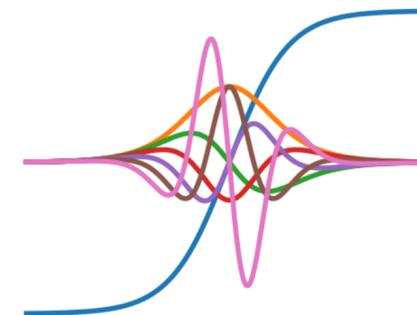
Advantages of automatic differentiation

- Accurate to the machine precision



- Same computational complexity as the function evaluation:
Baur-Strassen theorem '83

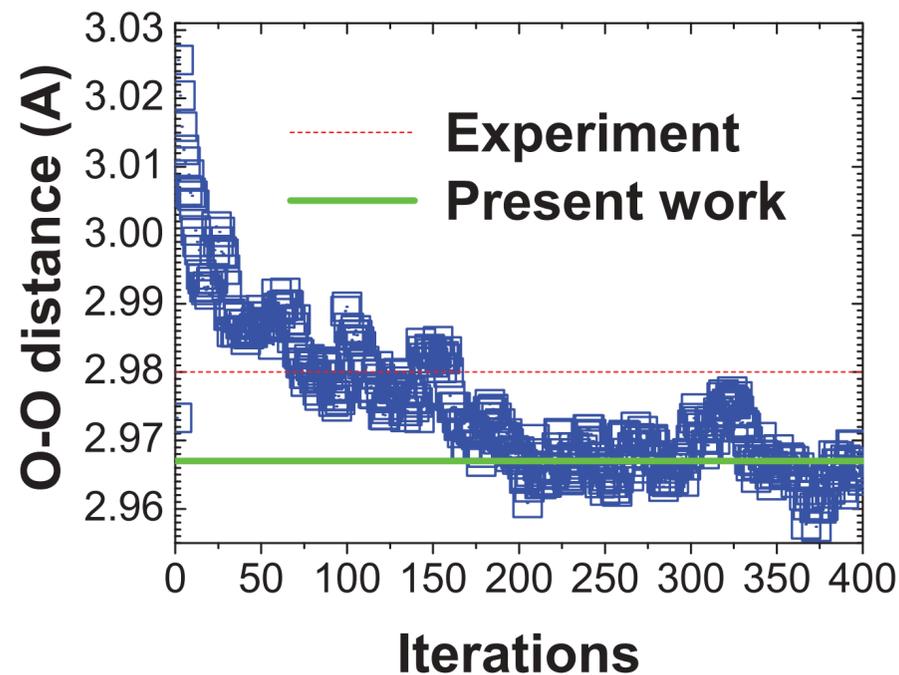
- Supports higher order gradients



```
>>> from autograd import elementwise_grad as egrad # for functions that vectorize over inputs
>>> import matplotlib.pyplot as plt
>>> x = np.linspace(-7, 7, 200)
>>> plt.plot(x, tanh(x),
...         x, egrad(tanh)(x), # first derivative
...         x, egrad(egrad(tanh))(x), # second derivative
...         x, egrad(egrad(egrad(tanh)))(x), # third derivative
...         x, egrad(egrad(egrad(egrad(tanh)))(x), # fourth derivative
...         x, egrad(egrad(egrad(egrad(egrad(tanh)))(x), # fifth derivative
...         x, egrad(egrad(egrad(egrad(egrad(egrad(tanh)))(x)) # sixth derivative
>>> plt.show()
```

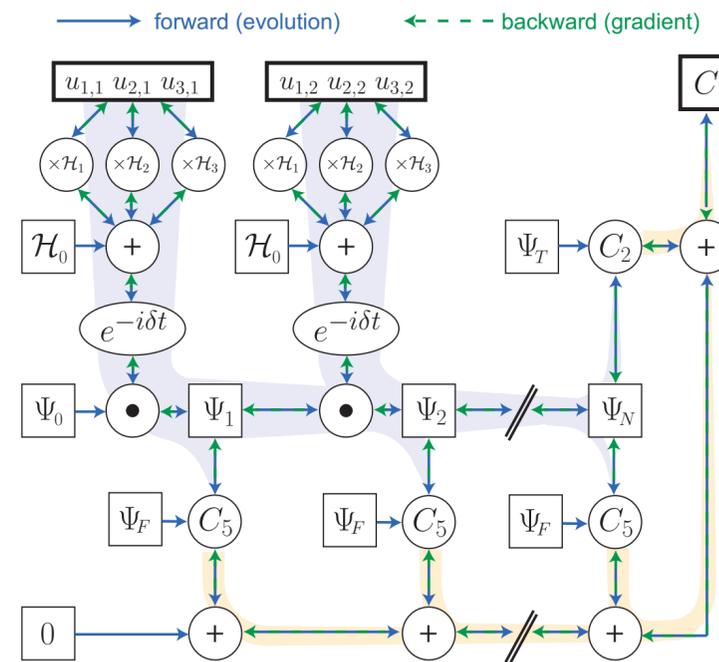
Applications of AD

Computing force



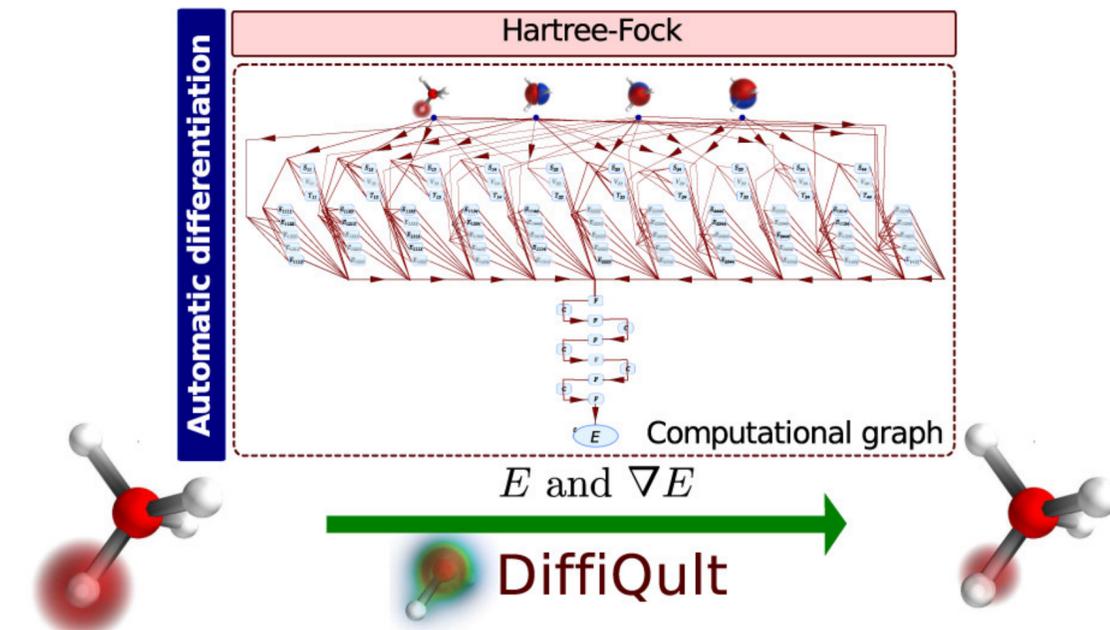
Sorella and Capriotti
J. Chem. Phys. '10

Quantum optimal control



Leung et al
PRA '17

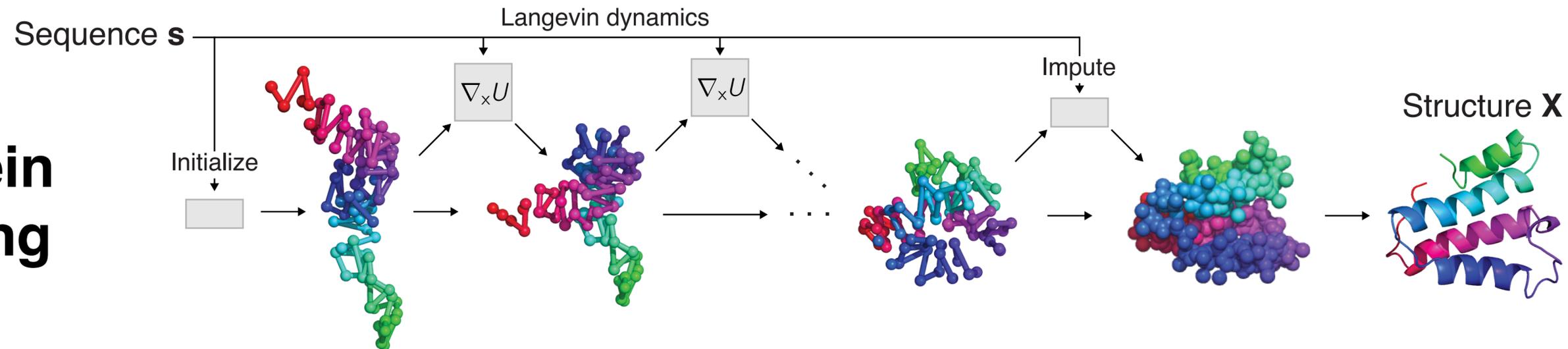
Variational Hartree-Fock



Tamayo-Mendoza et al
ACS Cent. Sci. '18

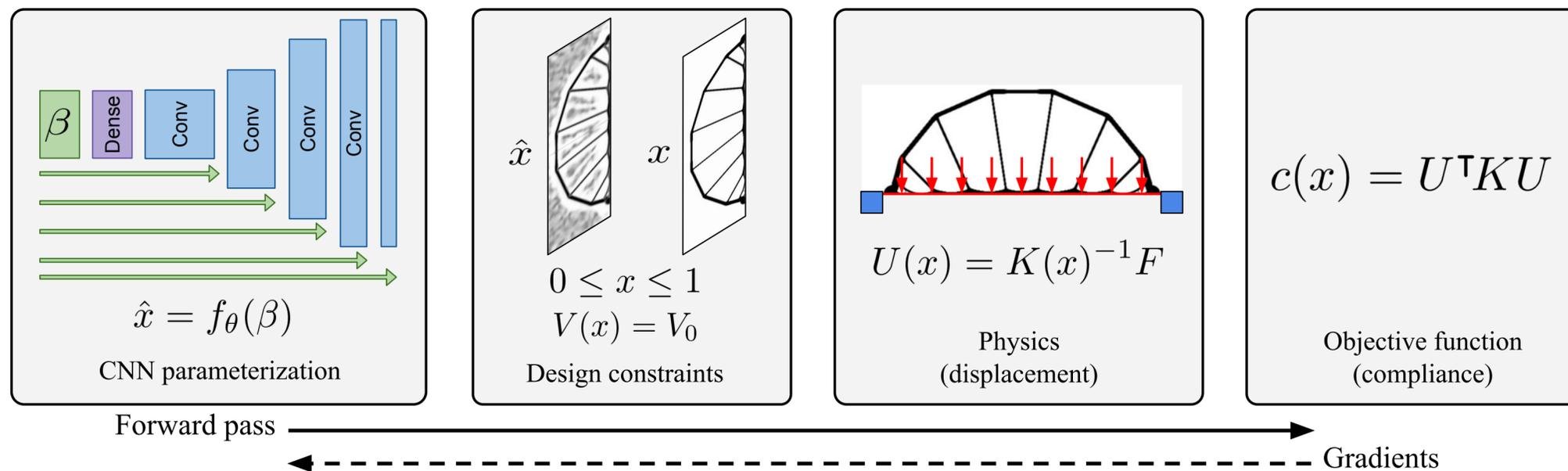
More Applications...

Protein folding

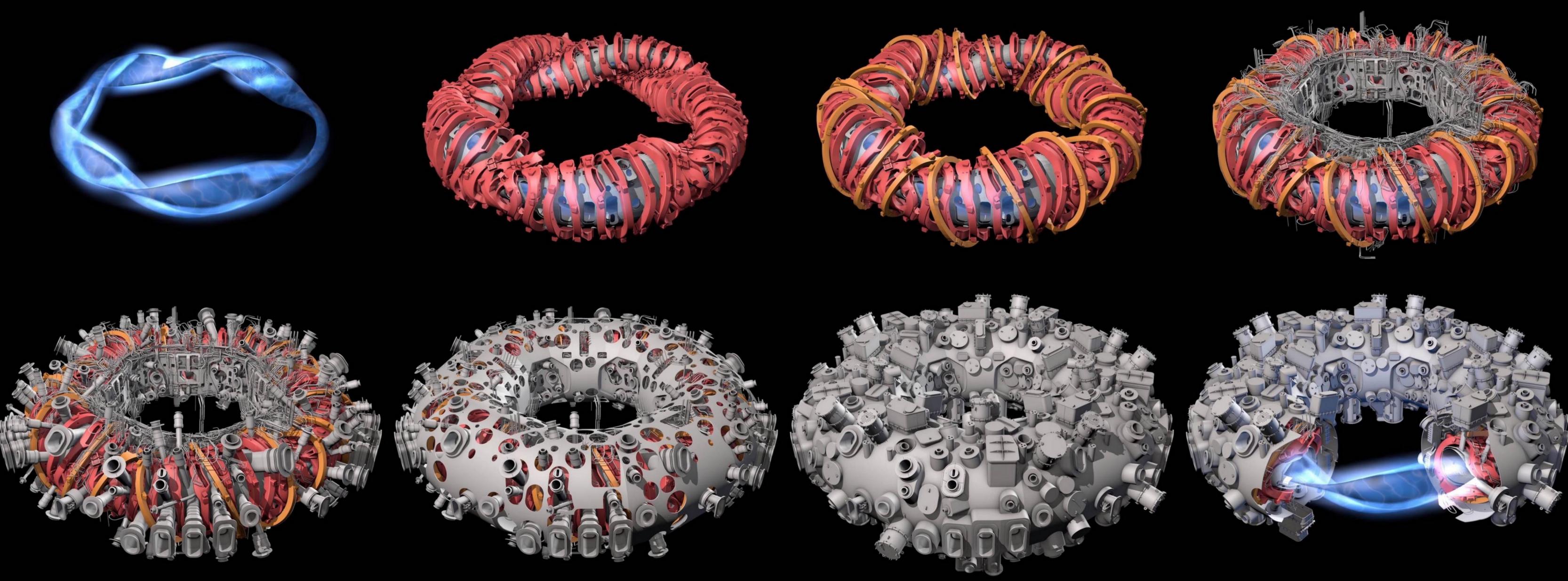


Ingraham et al
ICLR '19

Structural optimization



Hoyer et al
1909.04240



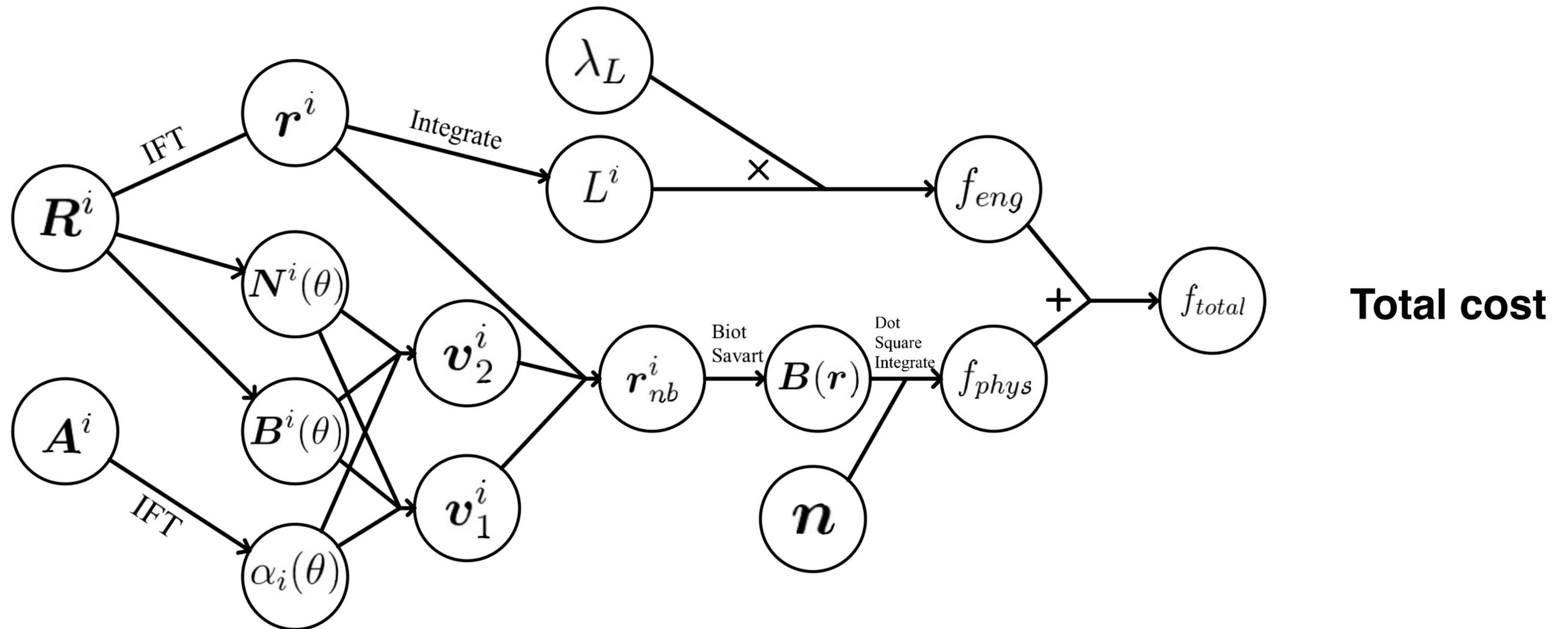
Coil design in fusion reactors (stellarator)

McGreivy et al 2009.00196

Computation graph

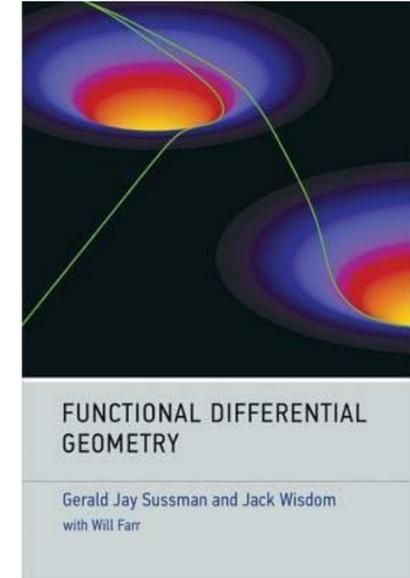
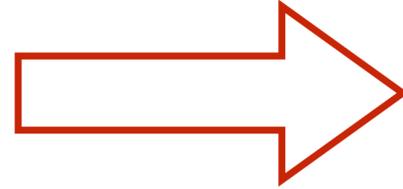
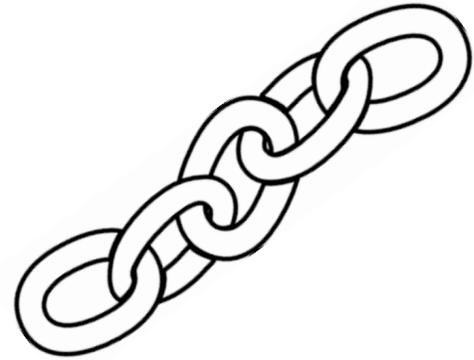
McGreivy et al 1909.04240

Coil parameters



Total cost

Differentiable programming is more than training neural networks



Black
magic box

Chain
rule



Functional
differential geometry

https://colab.research.google.com/github/google/jax/blob/master/notebooks/autodiff_cookbook.ipynb

Differentiating a general computer program (rather than neural networks) calls for deeper understanding of the technique

Reverse versus forward mode

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial x_n} \frac{\partial x_n}{\partial x_{n-1}} \dots \frac{\partial x_2}{\partial x_1} \frac{\partial x_1}{\partial \theta}$$


Reverse mode AD: **Vector-Jacobian Product of primitives**

- Backtrace the computation graph
- Needs to store intermediate results
- Efficient for graphs with large fan-in

Backpropagation = Reverse mode AD applied to neural networks

Reverse versus forward mode

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial x_n} \frac{\partial x_n}{\partial x_{n-1}} \dots \frac{\partial x_2}{\partial x_1} \frac{\partial x_1}{\partial \theta}$$

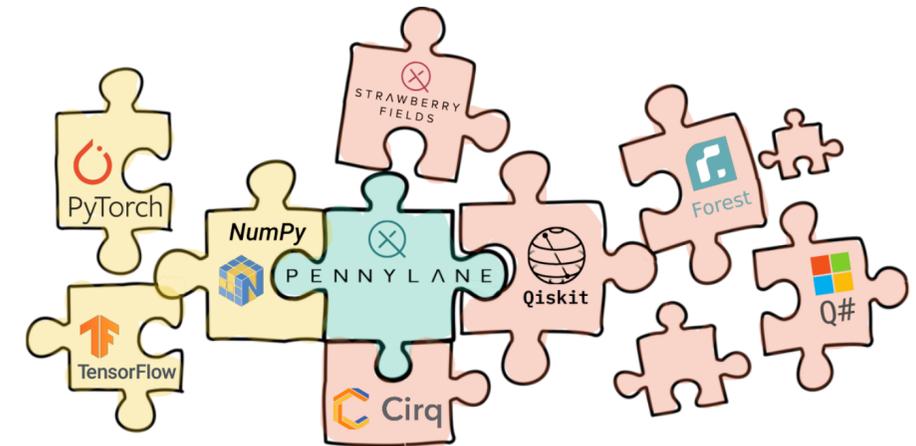
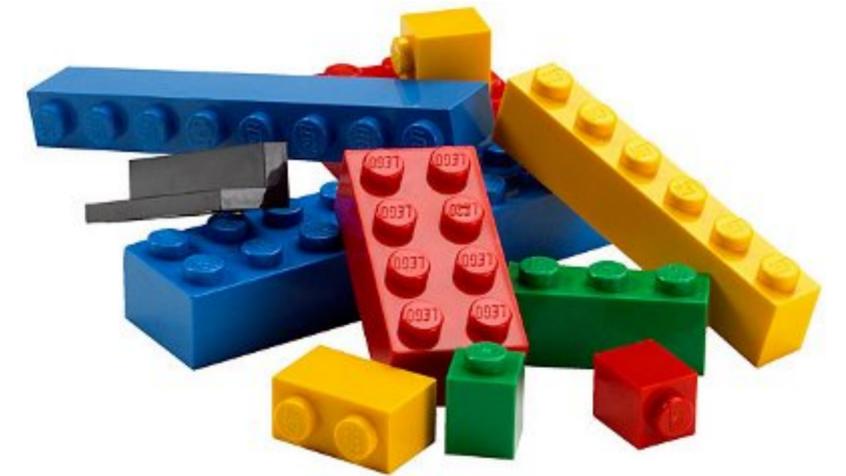

Forward mode AD: **Jacobian-Vector Product of primitives**

- Same order with the function evaluation
- **No storage overhead**
- Efficient for graph with large fan-out

Less efficient for scalar output, but useful for higher-order derivatives

How to think about AD ?

- AD is modular, and one can control its granularity
- Benefits of writing **customized primitives**
 - Reducing memory usage
 - Increasing numerical stability
- Call to **external libraries** written agnostically to AD
(or, even a quantum processor)



Example of primitives

~200 functions to cover most of numpy in HIPS/autograd

https://github.com/HIPS/autograd/blob/master/autograd/numpy/numpy_vjps.py

Operators	+, -, *, /, (-), **, %, <, <=, ==, !=, >=, >
Basic math functions	exp, log, square, sqrt, sin, cos, tan, sinh, cosh, tanh, sinc, abs, fabs, logaddexp, logaddexp2, absolute, reciprocal, exp2, expm1, log2, log10, log1p, arcsin, arccos, arctan, arcsinh, arccosh, arctanh, rad2deg, degrees, deg2rad, radians
Complex numbers	real, imag, conj, angle, fft, fftshift, ifftshift, real_if_close
Array reductions	sum, mean, prod, var, std, max, min, amax, amin
Array reshaping	reshape, ravel, squeeze, diag, roll, array_split, split, vsplit, hsplit, dsplit, expand_dims, flipud, fliplr, rot90, swapaxes, rollaxis, transpose, atleast_1d, atleast_2d, atleast_3d
Linear algebra	dot, tensordot, einsum, cross, trace, outer, det, slogdet, inv, norm, eigh, cholesky, sqrtm, solve_triangular
Other array operations	cumsum, clip, maximum, minimum, sort, msort, partition, concatenate, diagonal, truncate_pad, tile, full, triu, tril, where, diff, nan_to_num, vstack, hstack
Probability functions	t.pdf, t.cdf, t.logpdf, t.logcdf, multivariate_normal.logpdf, multivariate_normal.pdf, multivariate_normal.entropy, norm.pdf, norm.cdf, norm.logpdf, norm.logcdf,

```
67 # ----- Simple grads -----
68
69 defvjp(anp.negative, lambda ans, x: lambda g: -g)
70 defvjp(anp.abs,
71     lambda ans, x: lambda g: g * replace_zero(anp.conj(x), 0.) / replace_zero(ans, 1.))
72 defvjp(anp.fabs, lambda ans, x: lambda g: anp.sign(x) * g) # fabs doesn't take complex numbers.
73 defvjp(anp.absolute, lambda ans, x: lambda g: g * anp.conj(x) / ans)
74 defvjp(anp.reciprocal, lambda ans, x: lambda g: -g / x**2)
75 defvjp(anp.exp, lambda ans, x: lambda g: ans * g)
76 defvjp(anp.exp2, lambda ans, x: lambda g: ans * anp.log(2) * g)
77 defvjp(anp.expm1, lambda ans, x: lambda g: (ans + 1) * g)
78 defvjp(anp.log, lambda ans, x: lambda g: g / x)
79 defvjp(anp.log2, lambda ans, x: lambda g: g / x / anp.log(2))
80 defvjp(anp.log10, lambda ans, x: lambda g: g / x / anp.log(10))
81 defvjp(anp.log1p, lambda ans, x: lambda g: g / (x + 1))
82 defvjp(anp.sin, lambda ans, x: lambda g: g * anp.cos(x))
83 defvjp(anp.cos, lambda ans, x: lambda g: -g * anp.sin(x))
84 defvjp(anp.tan, lambda ans, x: lambda g: g / anp.cos(x)**2)
85 defvjp(anp.arcsin, lambda ans, x: lambda g: g / anp.sqrt(1 - x**2))
86 defvjp(anp.arccos, lambda ans, x: lambda g: -g / anp.sqrt(1 - x**2))
87 defvjp(anp.arctan, lambda ans, x: lambda g: g / (1 + x**2))
88 defvjp(anp.sinh, lambda ans, x: lambda g: g * anp.cosh(x))
89 defvjp(anp.cosh, lambda ans, x: lambda g: g * anp.sinh(x))
90 defvjp(anp.tanh, lambda ans, x: lambda g: g / anp.cosh(x)**2)
91 defvjp(anp.arcsinh, lambda ans, x: lambda g: g / anp.sqrt(x**2 + 1))
92 defvjp(anp.arccosh, lambda ans, x: lambda g: g / anp.sqrt(x**2 - 1))
93 defvjp(anp.arctanh, lambda ans, x: lambda g: g / (1 - x**2))
94 defvjp(anp.rad2deg, lambda ans, x: lambda g: g / anp.pi * 180.0)
95 defvjp(anp.degrees, lambda ans, x: lambda g: g / anp.pi * 180.0)
96 defvjp(anp.deg2rad, lambda ans, x: lambda g: g * anp.pi / 180.0)
97 defvjp(anp.radians, lambda ans, x: lambda g: g * anp.pi / 180.0)
98 defvjp(anp.square, lambda ans, x: lambda g: g * 2 * x)
99 defvjp(anp.sqrt, lambda ans, x: lambda g: g * 0.5 * x**-0.5)
```

Loop/Condition/Sort/Permutations are also differentiable

Differentiable programming tools

HIPS/autograd

theano



 PyTorch


TensorFlow



 Keras



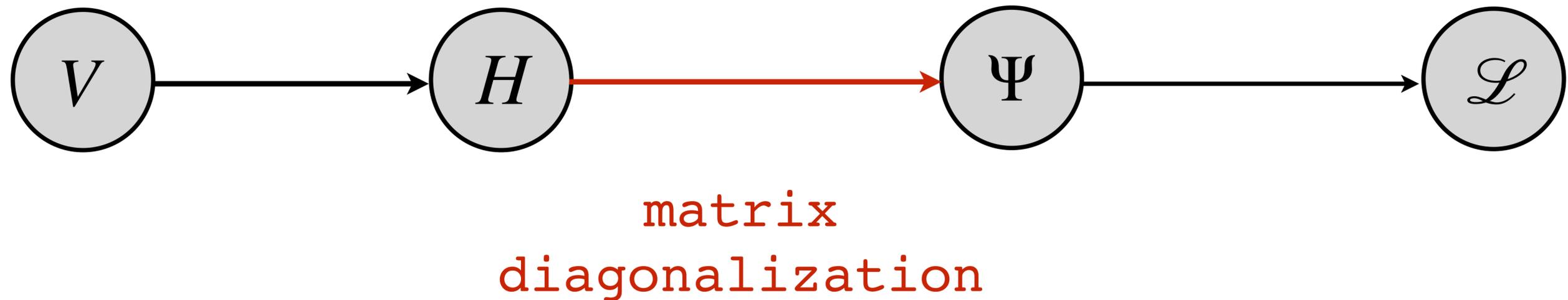
Differentiable Scientific Computing

- Many scientific computations (**FFT, Eigen, SVD!**) are [differentiable](#)
- ODE integrators are differentiable with [O\(1\) memory](#)
- [Differentiable ray tracer](#) and [Differentiable fluid simulations](#)
- Differentiable Monte Carlo/Tensor Network/Functional RG/
Dynamical Mean Field Theory/Density Functional Theory/
Hartree-Fock/Coupled Cluster/Gutzwiller/Molecular Dynamics...

Differentiate through domain-specific computational processes to solve learning, control, optimization and inverse problems

Differentiable Eigensolver

Inverse Schrodinger Problem



Differentiable Eigensolver

$$H\Psi = \Psi E$$

Forward mode: What happen if $H \rightarrow H + dH$? Perturbation theory

Reverse mode: How should I change H given $\partial\mathcal{L}/\partial\Psi$ and $\partial\mathcal{L}/\partial E$? **Inverse perturbation theory!**

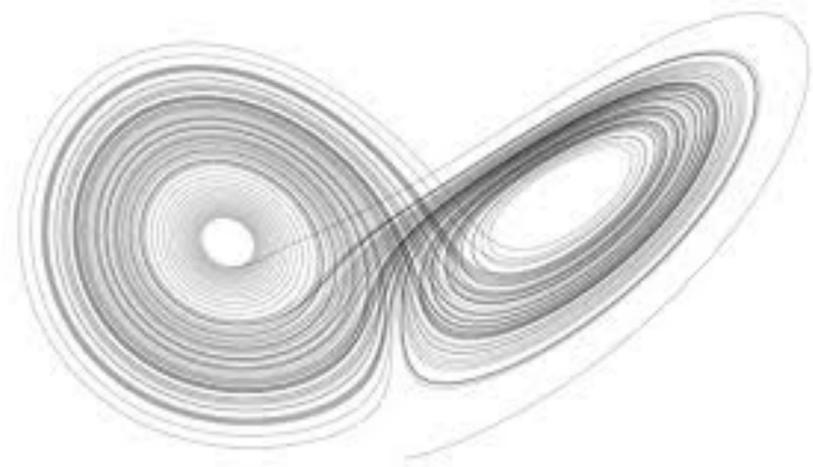
Hamiltonian engineering via differentiable programming



Differentiable ODE integrators

“Neural ODE” Chen et al, 1806.07366

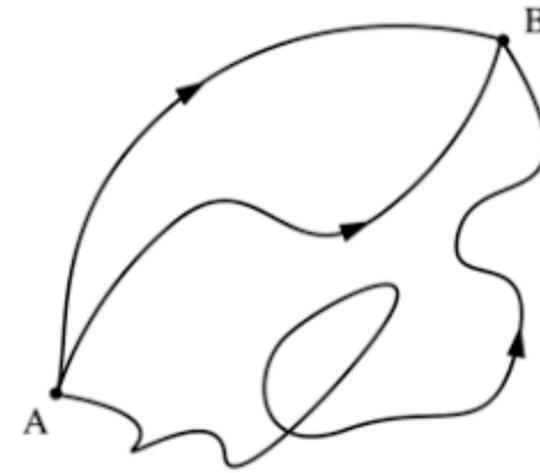
Dynamics systems



$$\frac{dx}{dt} = f_{\theta}(x, t)$$

Classical and quantum control

Principle of least actions

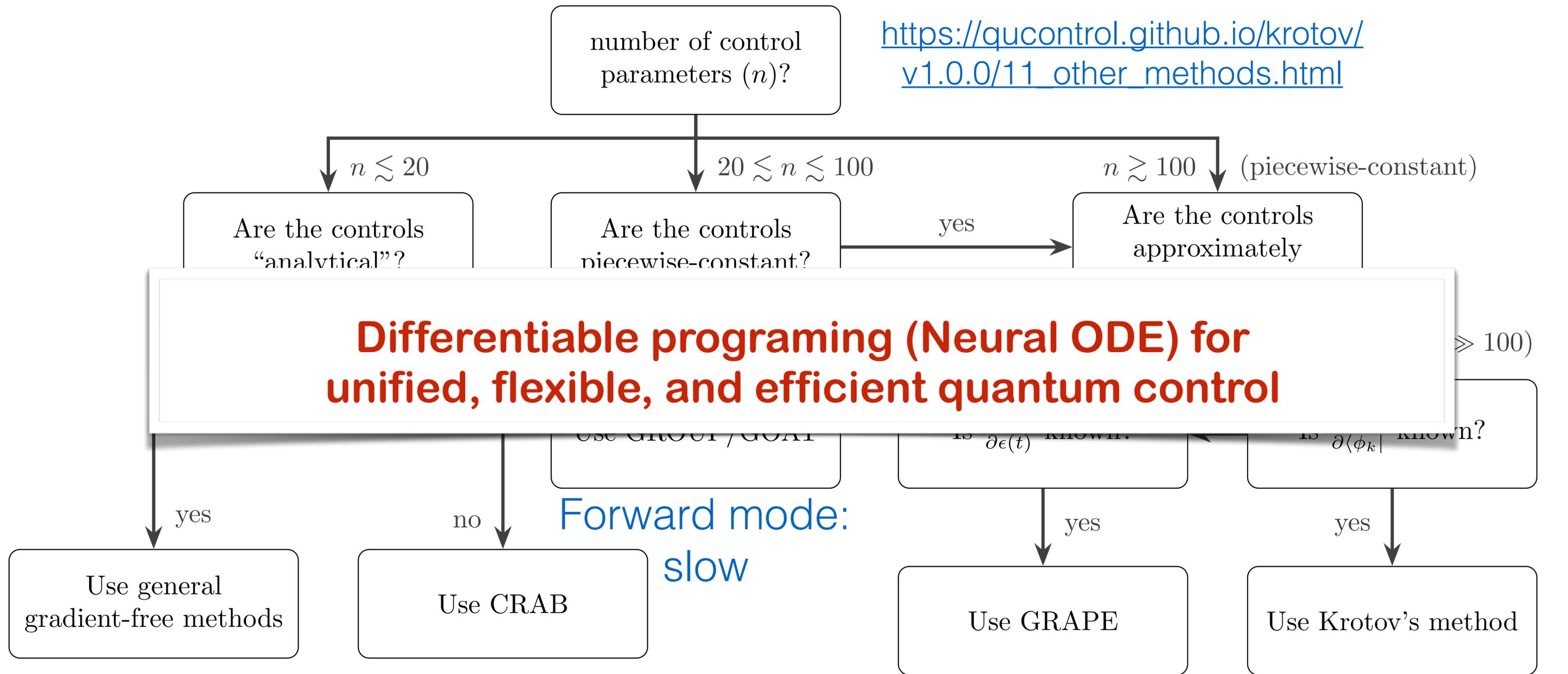


$$S = \int \mathcal{L}(q_{\theta}, \dot{q}_{\theta}, t) dt$$

Optics, (quantum) mechanics, field theory...

Quantum optimal control $i\frac{dU}{dt} = HU$

https://qucontrol.github.io/krotov/v1.0.0/11_other_methods.html



No gradient:
not scalable

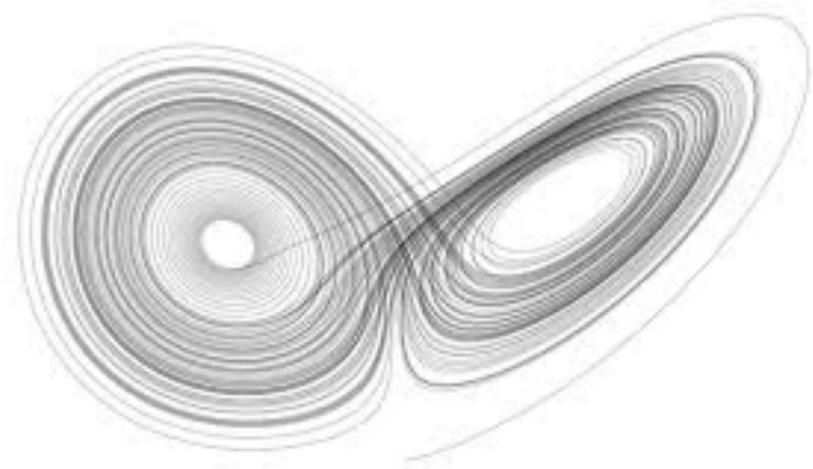
Forward mode:
slow

Reverse mode w/ discretize steps:
piecewise-constant assumption

Differentiable ODE integrators

“Neural ODE” Chen et al, 1806.07366

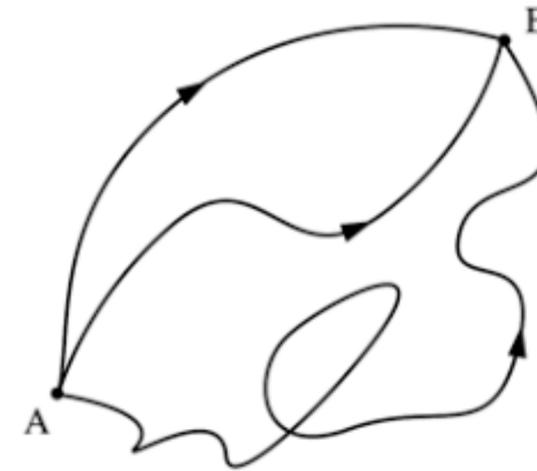
Dynamics systems



$$\frac{dx}{dt} = f_{\theta}(x, t)$$

Classical and quantum control

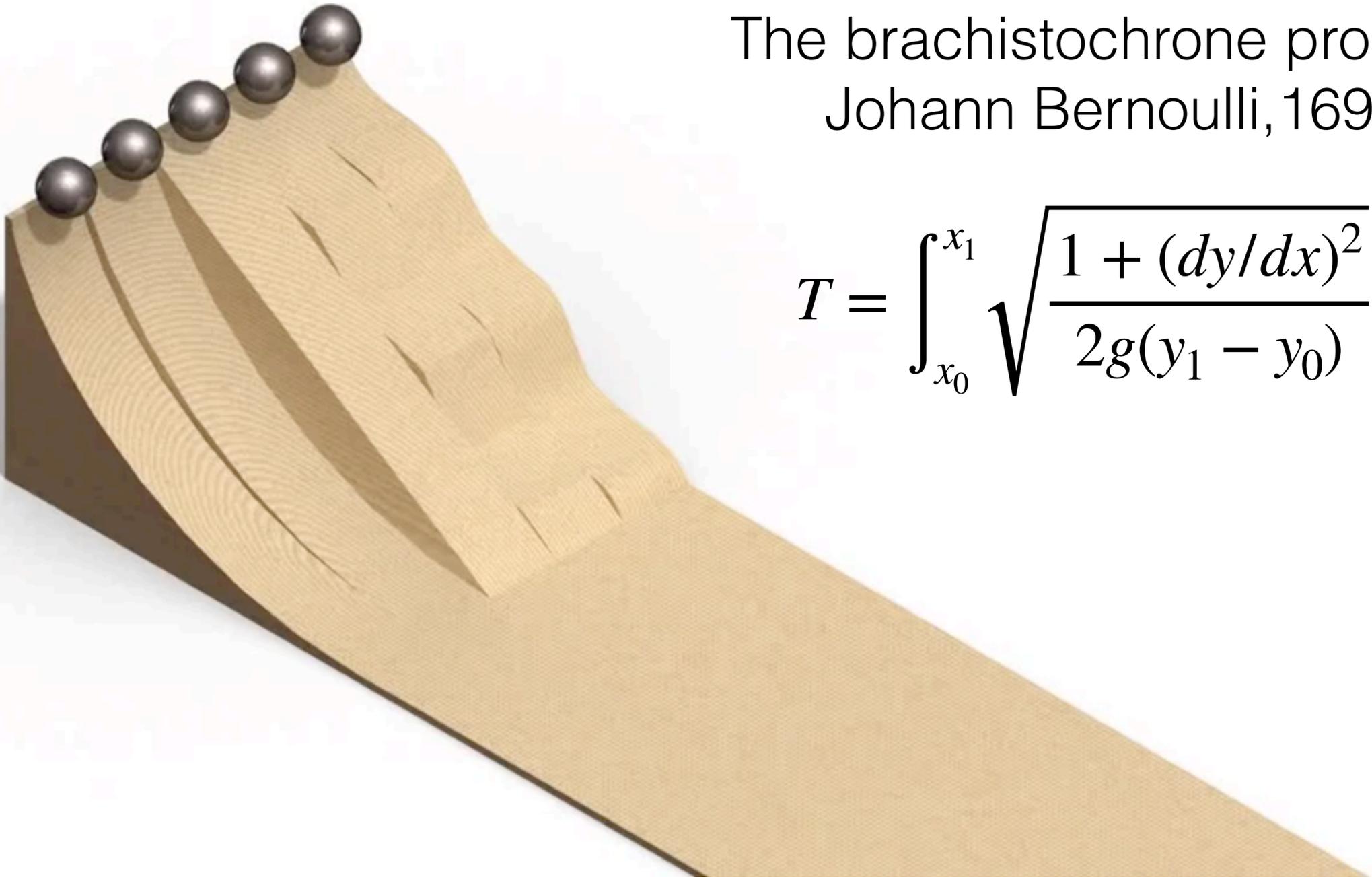
Principle of least actions



$$S = \int \mathcal{L}(q_{\theta}, \dot{q}_{\theta}, t) dt$$

Optics, (quantum) mechanics, field theory...

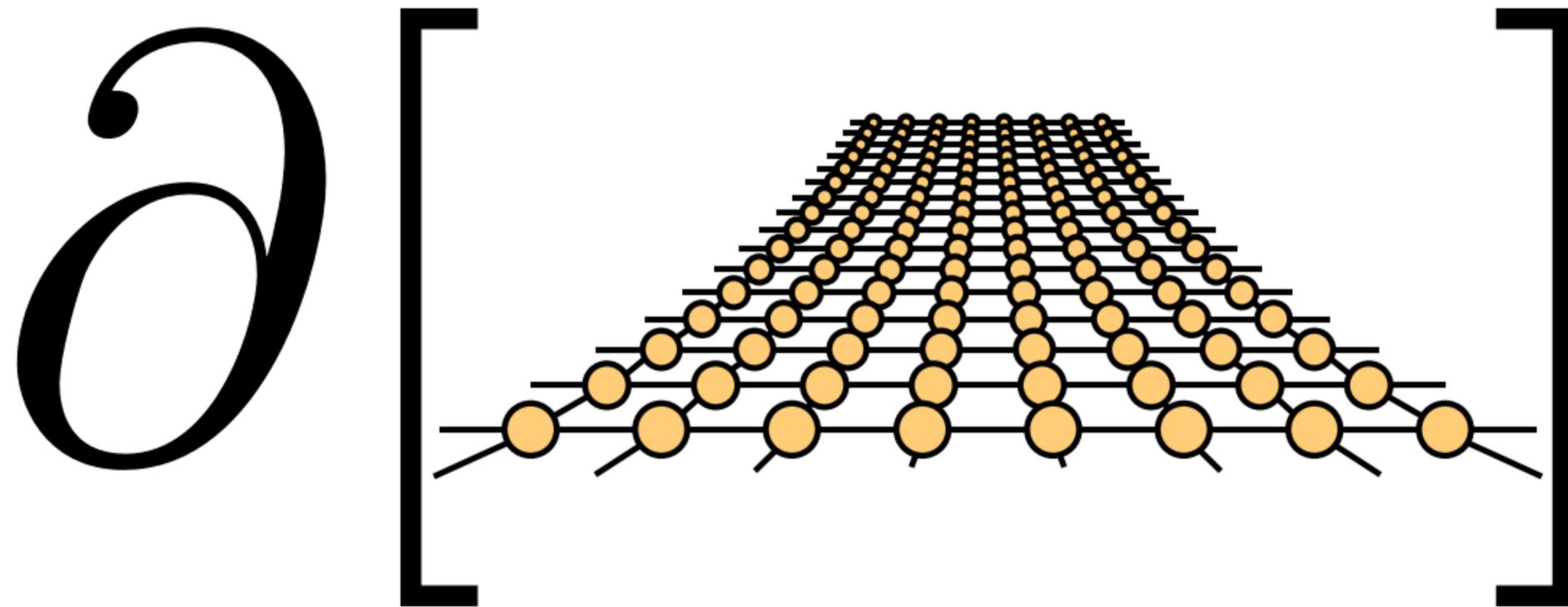
Differentiable functional optimization



The brachistochrone problem
Johann Bernoulli, 1696

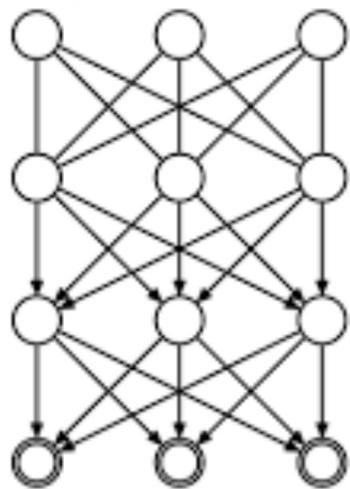
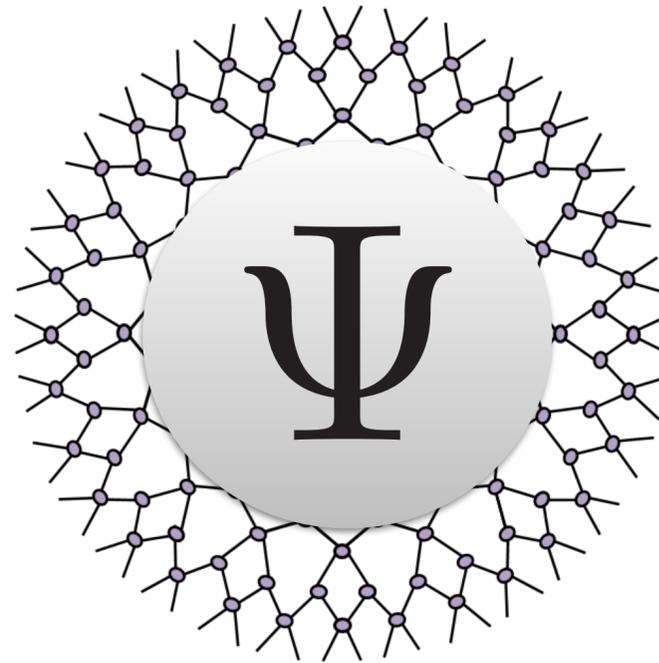
$$T = \int_{x_0}^{x_1} \sqrt{\frac{1 + (dy/dx)^2}{2g(y_1 - y_0)}} dx$$

Differentiable Programming Tensor Networks

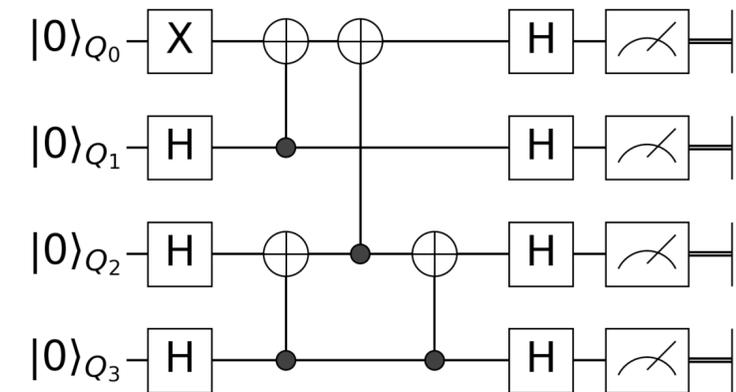


“Tensor network is 21 century’s matrix”

—Mario Szegedy



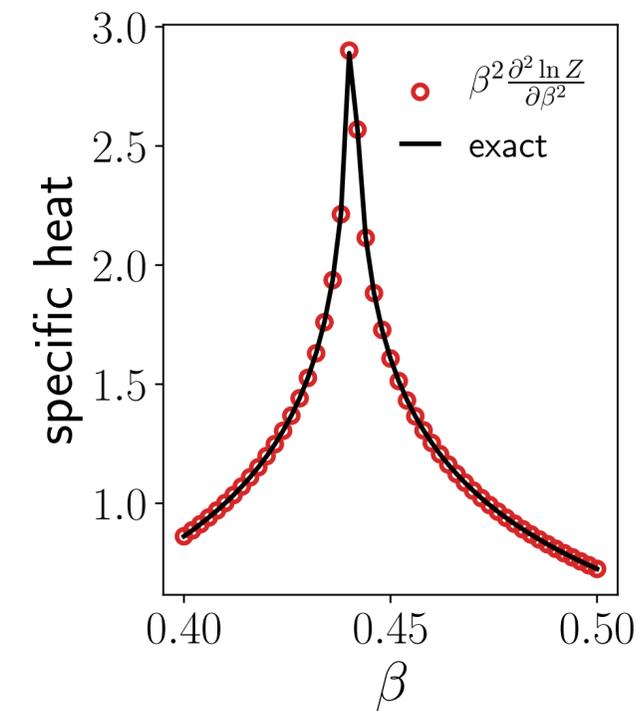
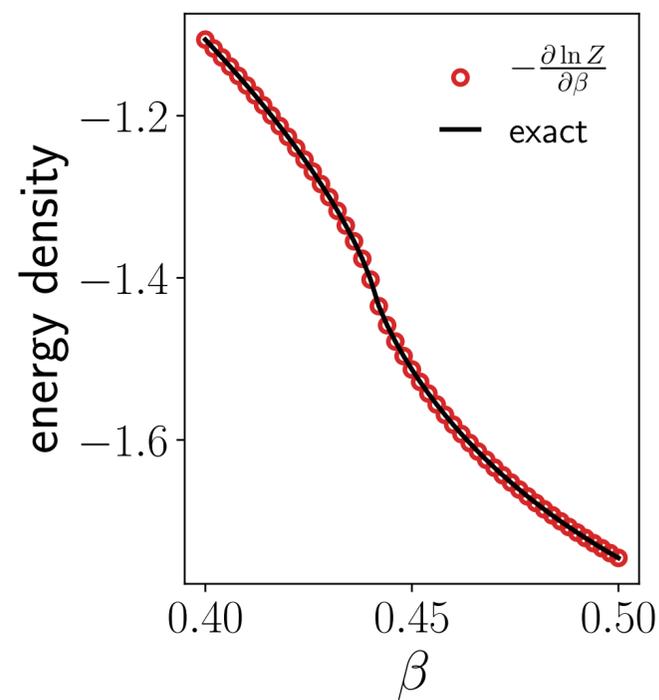
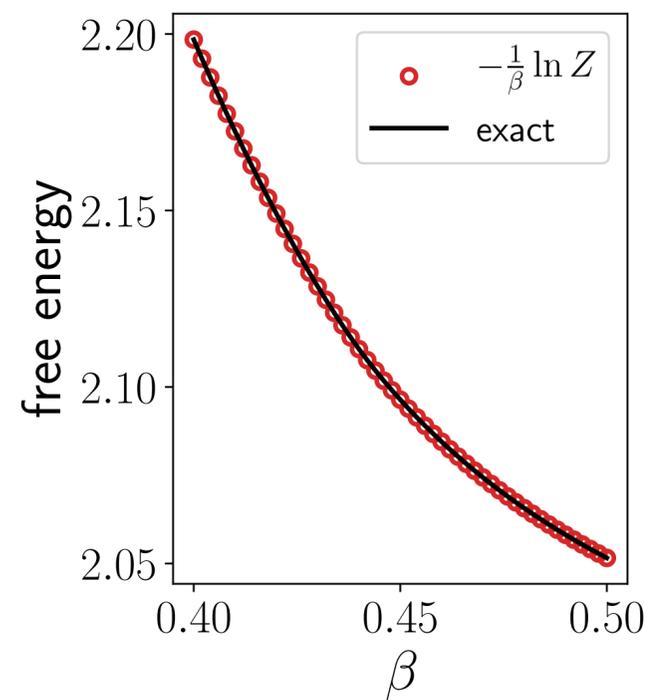
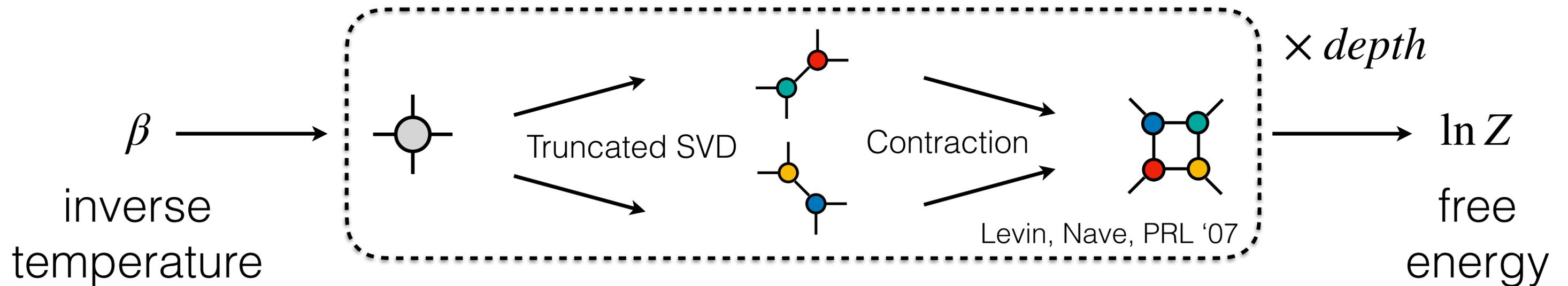
**Neural networks and
Probabilistic graphical models**



**Quantum circuit architecture,
parametrization, and simulation**

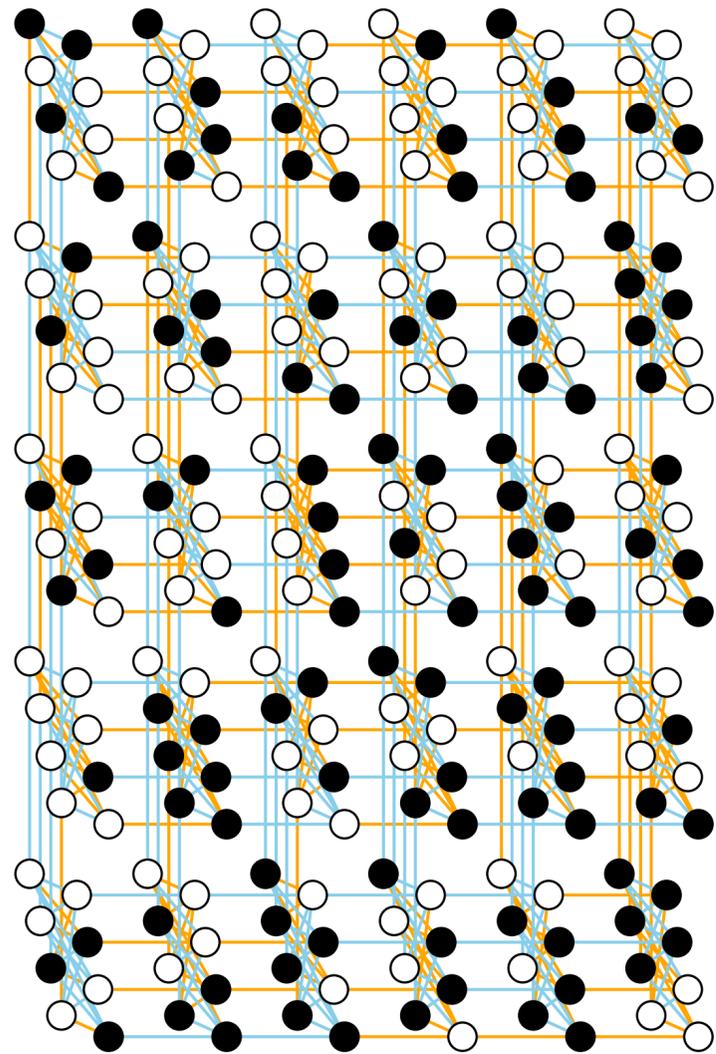
Differentiate through tensor renormalization group

Computation graph



Compute physical observables as gradient of tensor network contraction

Differentiable spin glass solver



couplings
& fields

tensor network
contraction

optimal
energy

optimal
configuration

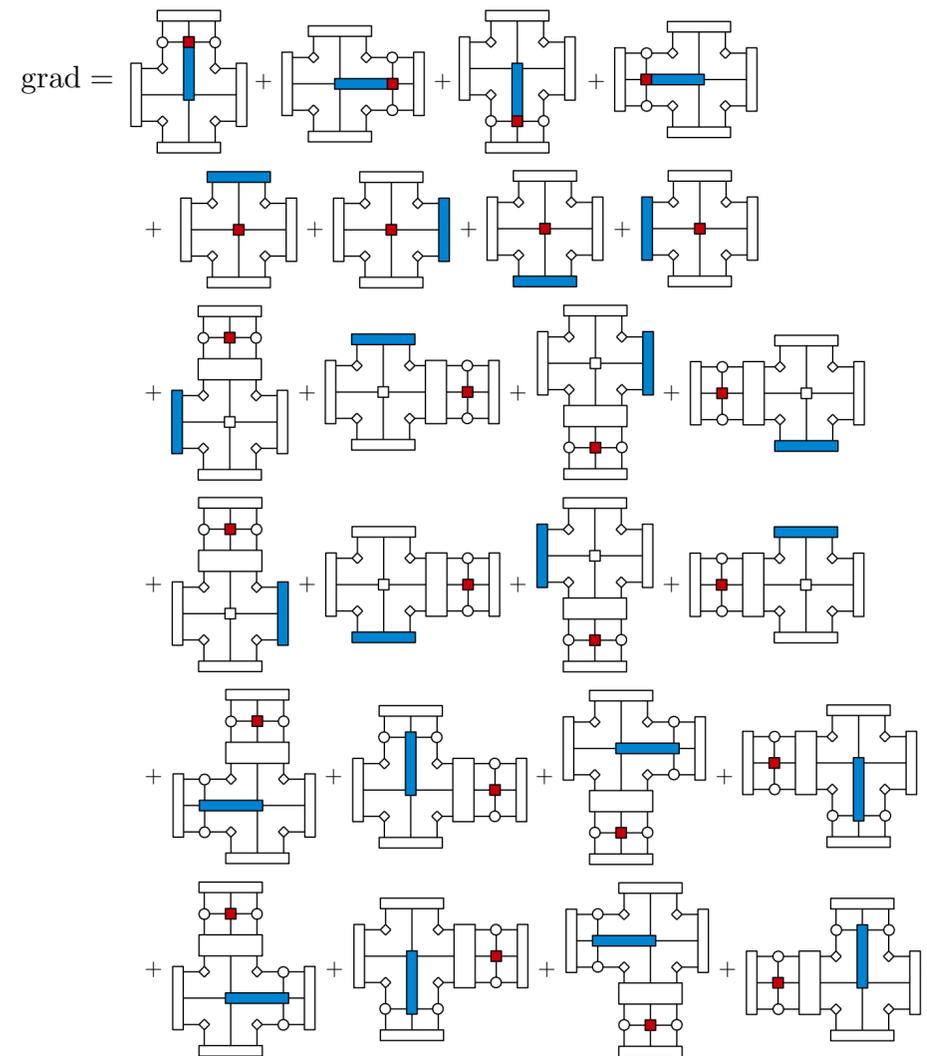
=

$$\frac{\partial [\text{optimal energy}]}{\partial [\text{field}]}$$



Differentiable iPEPS optimization

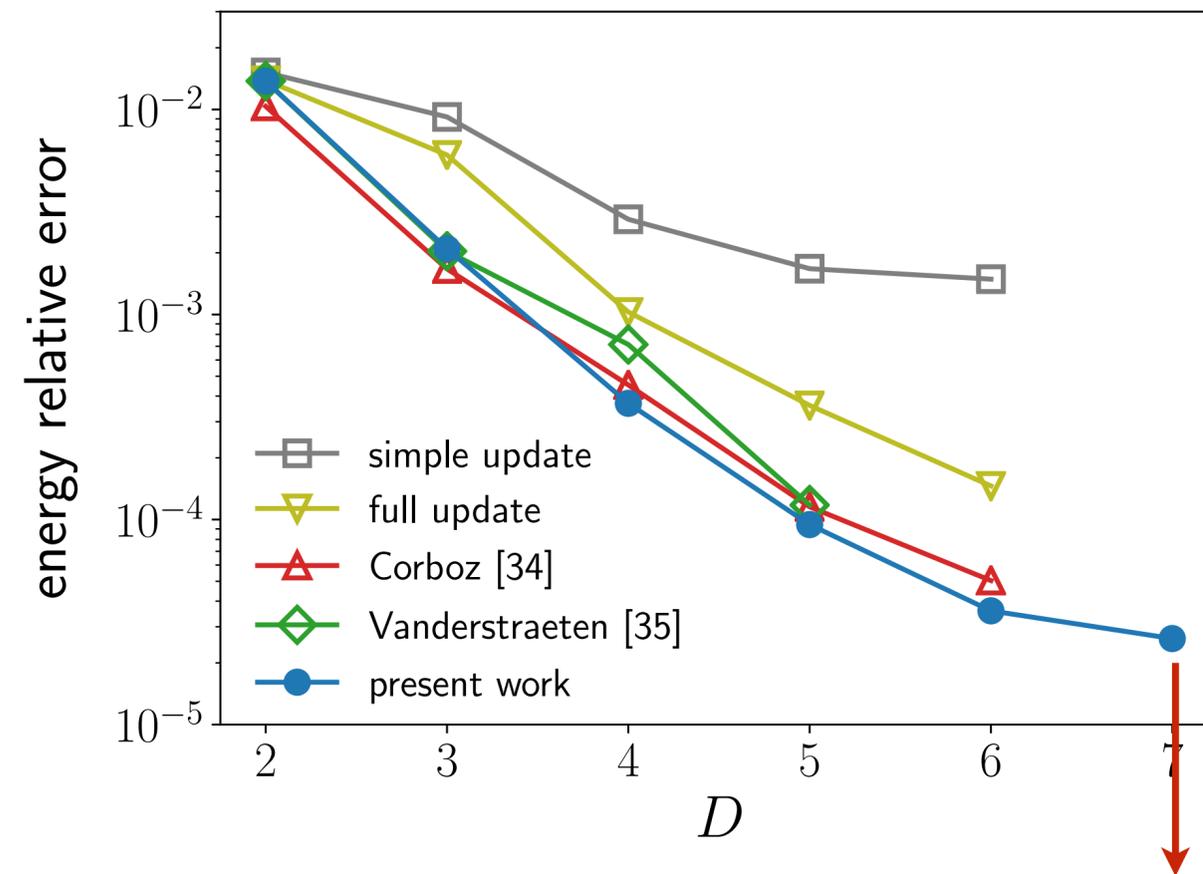
before...



Vanderstraeten et al, PRB '16

now, w/ differentiable programming

Liao, Liu, LW, Xiang, PRX '19



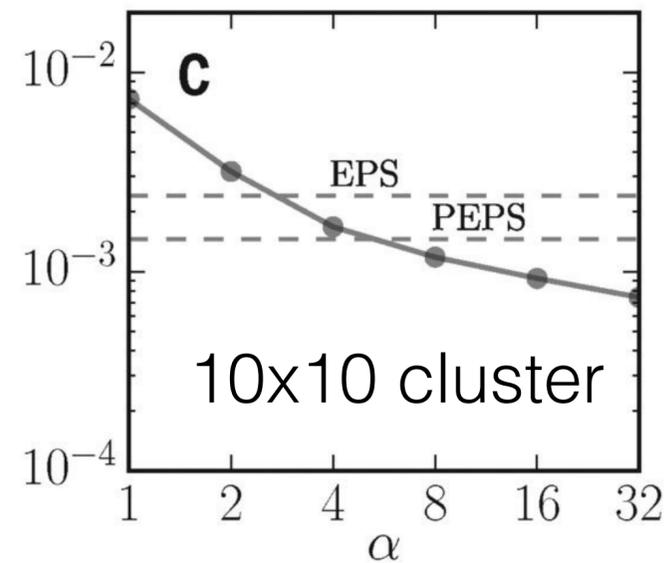
Best variational energy to date

<https://github.com/wangleiphy/tensorgrad> 1 GPU (Nvidia P100) week

Differentiable iPEPS optimization

Finite size

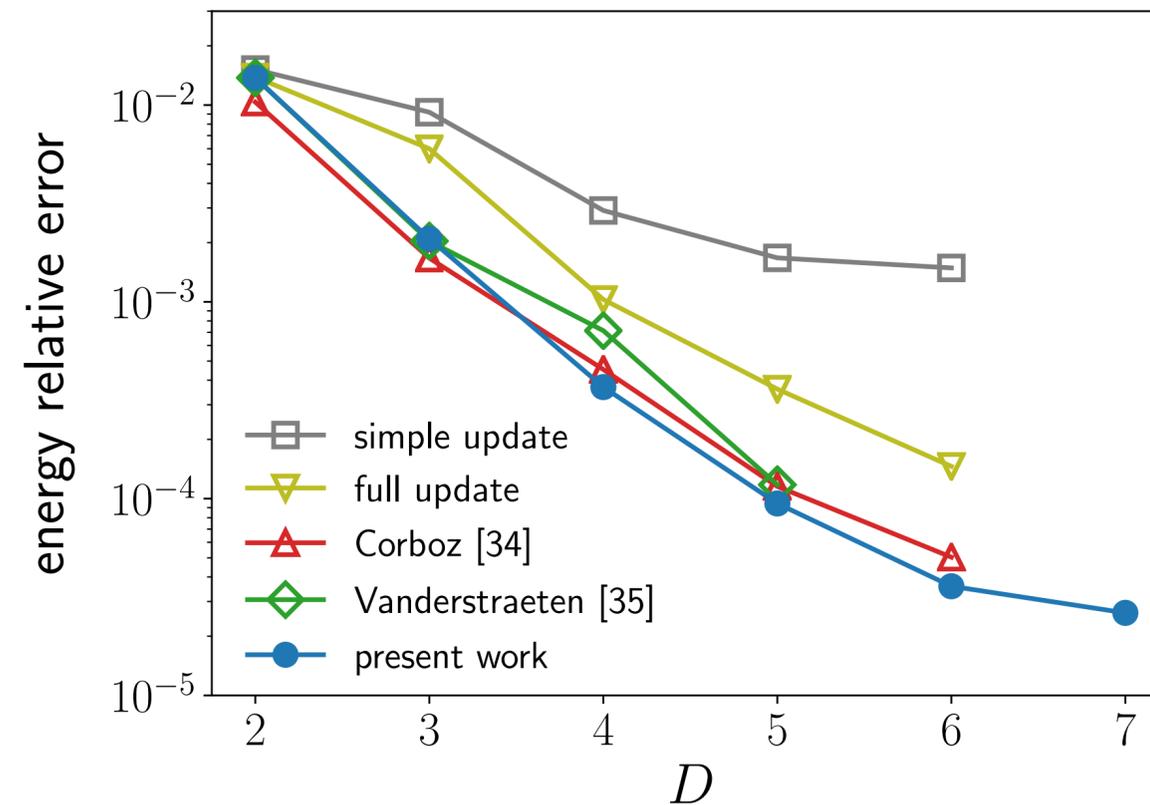
Neural network



Carleo & Troyer, Science '17

Infinite size

Tensor network



Liao, Liu, LW, Xiang, PRX '19

**Further progress for challenging physical problems:
frustrated magnets, fermions, thermodynamics ...**

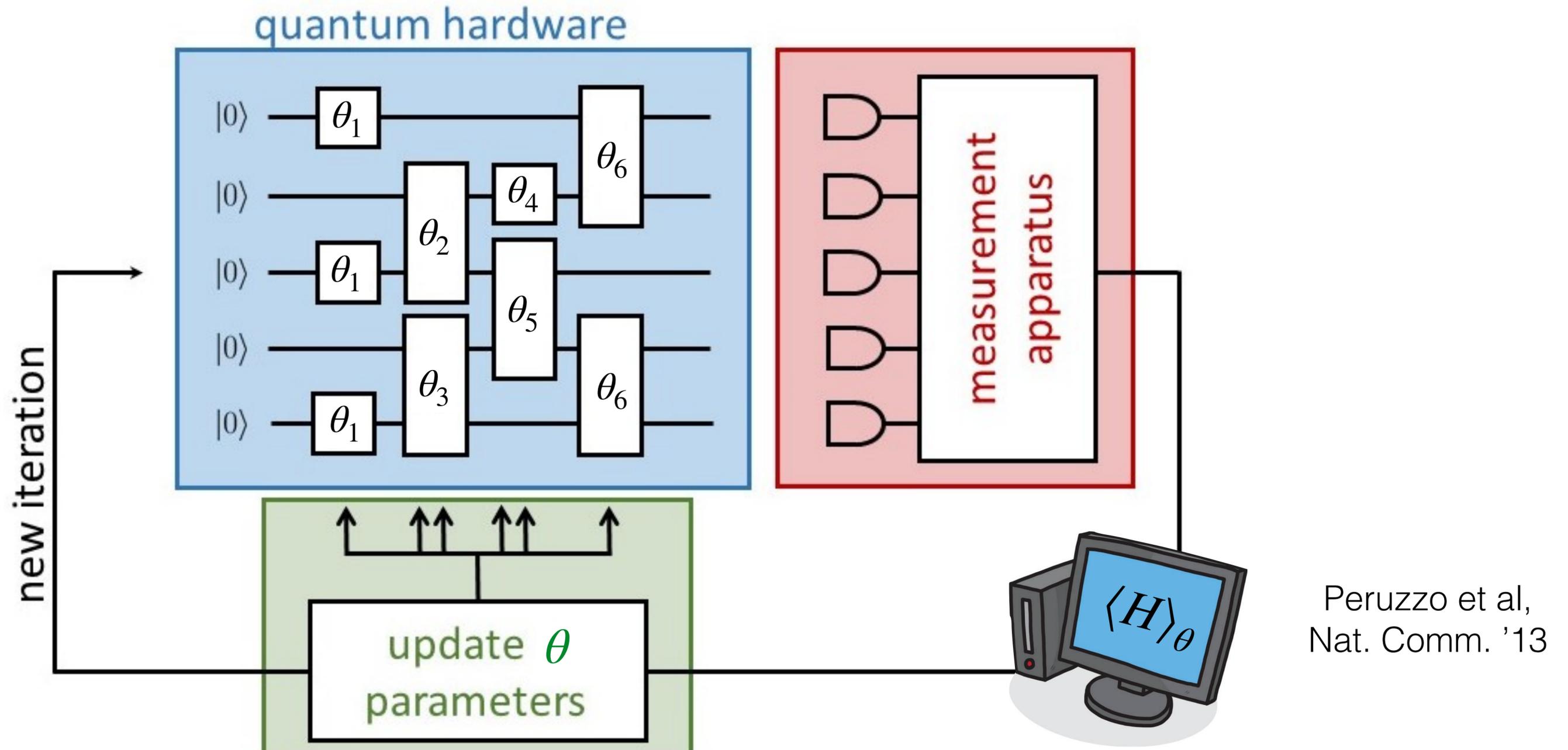
Chen et al, '19
Xie et al, '20
Tang et al '20

Differentiable Programming

Quantum Circuits

neural networks — graphical models — tensor networks — quantum circuits

Variational quantum algorithms

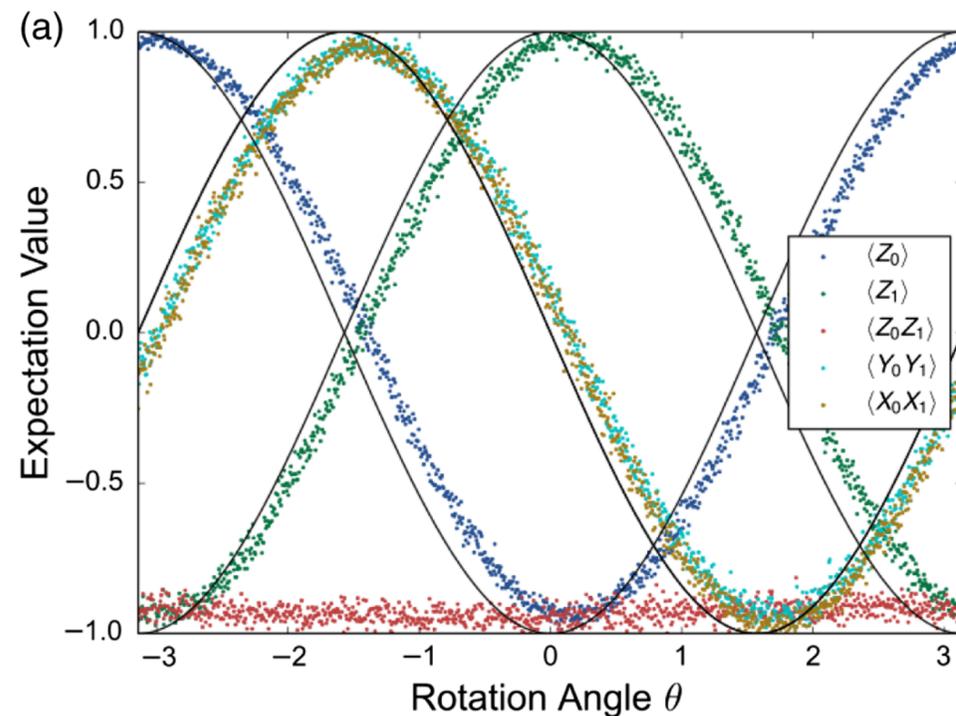


Peruzzo et al,
Nat. Comm. '13

Quantum circuit as a variational ansatz

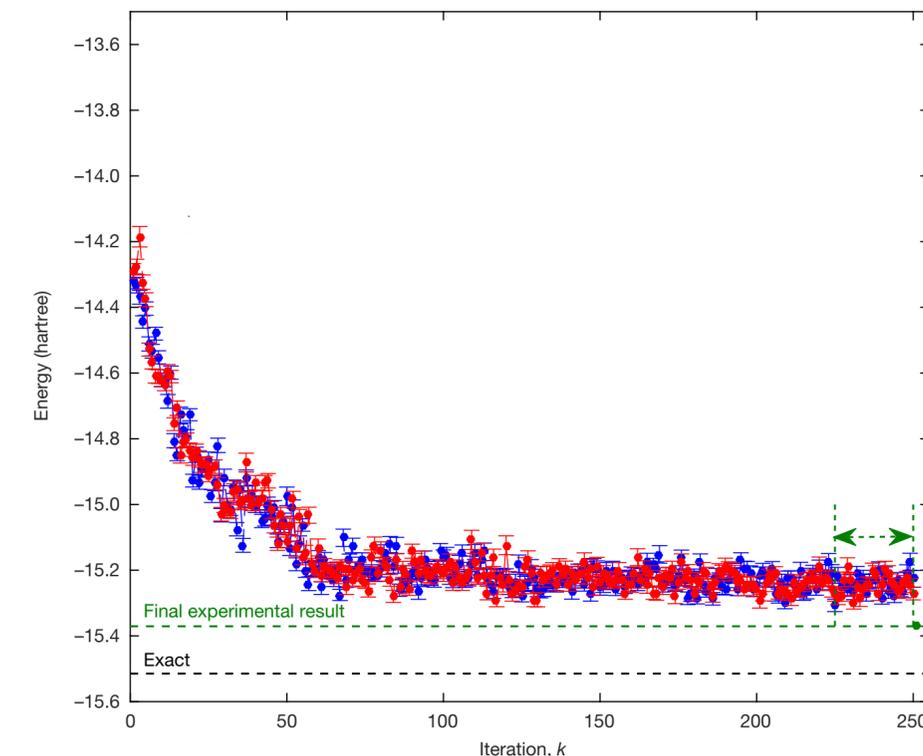
Optimize variational quantum circuits

Scan the single variational parameter



Google PRX '16

Stochastic perturbation of 30 variational parameters

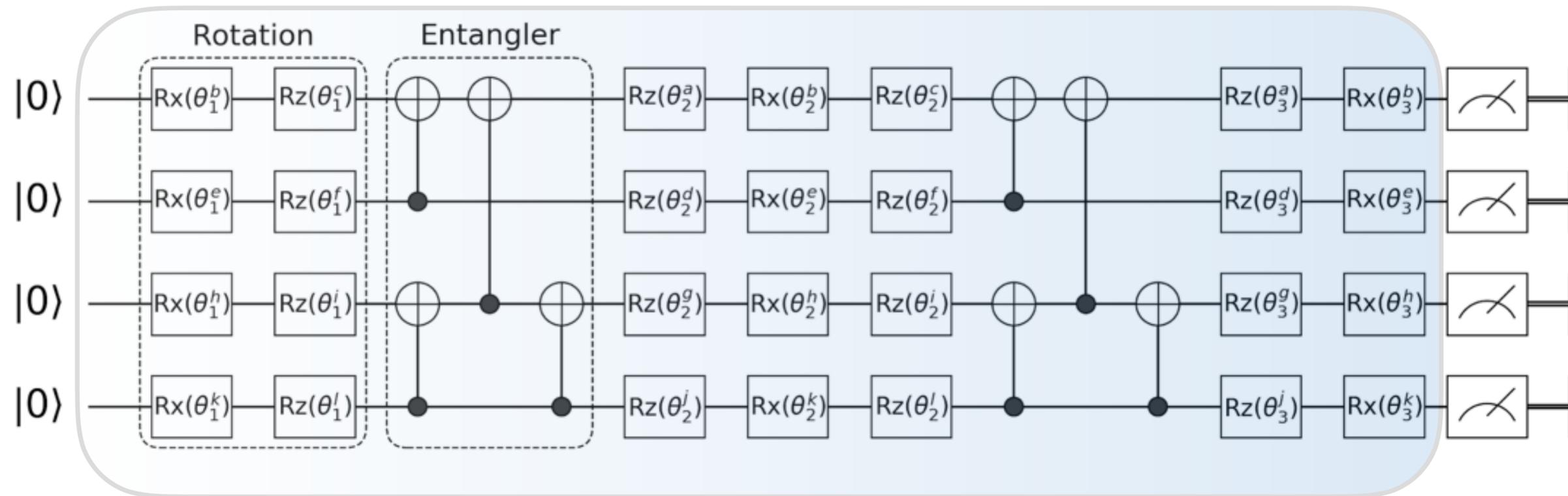


IBM Nature '17

Optimization with analytical gradient is essential for higher dimensions

Differentiable¹ quantum circuits

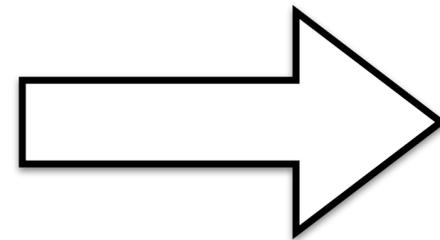
measure gradient on real device



Parametrized gate of the form

$$e^{-\frac{i\theta}{2}\Sigma} \text{ with } \Sigma^2 = 1$$

e.g., X, Y, Z, CNOT, SWAP...



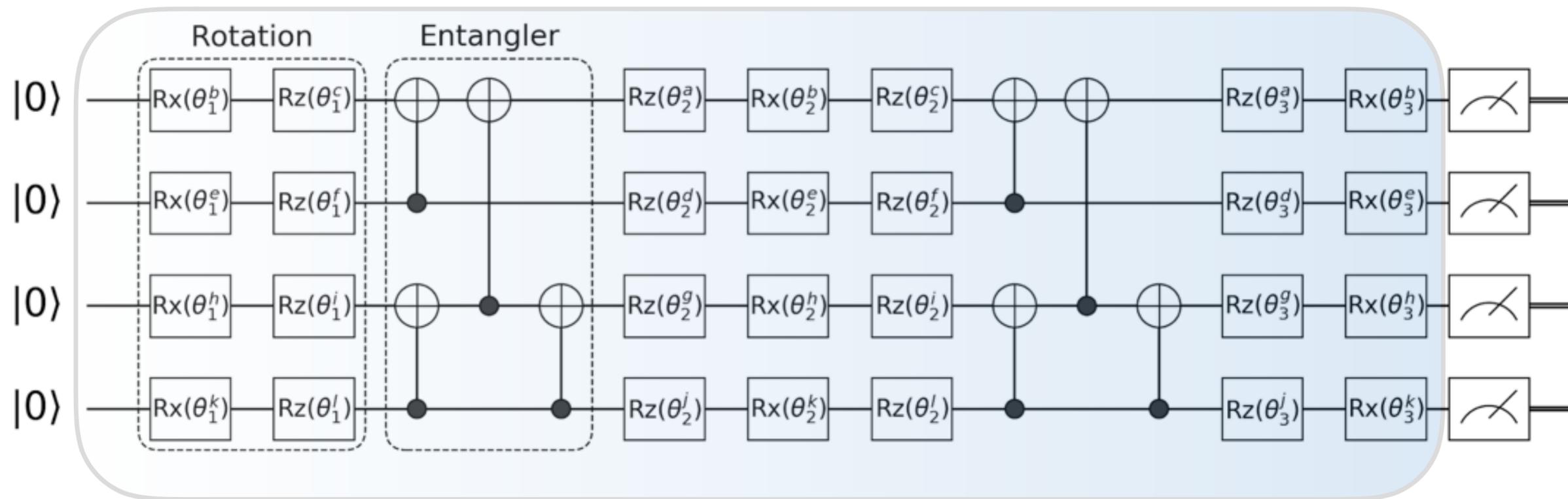
Li et al, PRL '17, Mitarai et al, PRA '18
Schuld et al, PRA '19, Crooks, '19...

$$\nabla \langle H \rangle_\theta = \left(\langle H \rangle_{\theta+\pi/2} - \langle H \rangle_{\theta-\pi/2} \right) / 2$$

Same complexity as forward mode automatic differentiation

*Differentiable*² quantum circuits

compute gradient in classical simulations



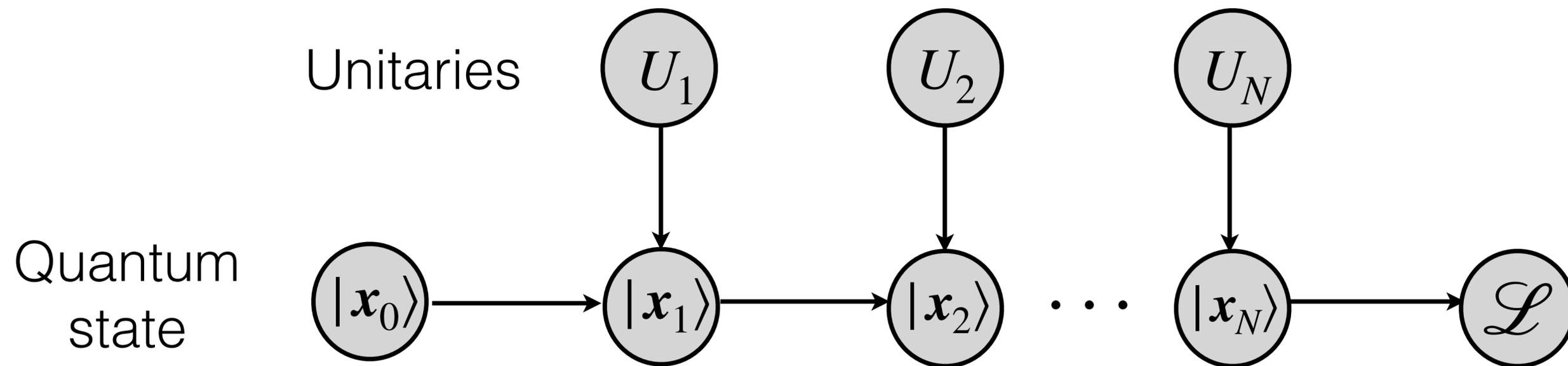
P E N N Y L A N E



TensorFlow Quantum

Unfortunately, forward mode is slow
Reverse mode is memory consuming

Quantum circuit computation graph



The same “comb graph” as the feedforward neural network, except that quantum computing is reversible

Reversible AD for variational quantum circuits*

forward

$$U|x\rangle \rightarrow |y\rangle$$

backward

“uncompute” $|x\rangle \leftarrow U^\dagger |y\rangle$

adjoint
for mat-vec
multiply

$$\overline{|x\rangle} \leftarrow U^\dagger \overline{|y\rangle}$$

$$\overline{U} \leftarrow \overline{|y\rangle} \langle x|$$

All are in-place operations without caching

*GRAPE type algorithm on the level of circuits

```

julia> using Yao, YaoExtensions

julia> n = 10; depth = 10000;

julia> circuit = dispatch!(
    variational_circuit(n, depth),
    :random);

julia> gatecount(circuit)
Dict{Type{#s54} where #s54 <:
    AbstractBlock,Int64} with 3 entries:
  RotationGate{1,Float64,ZGate} => 200000
  RotationGate{1,Float64,XGate} => 100010
  ControlBlock{10,XGate,1,1}    => 100000

julia> nparameters(circuit)
300010

julia> h = heisenberg(n);

julia> for i = 1:100
    _, grad = expect'(h, zero_state(n)=>
                      circuit)
    dispatch!(-, circuit, 1e-3 * grad)
    println("Step $i, energy = $(expect(
        h, zero_state(n)=>circuit))")
end

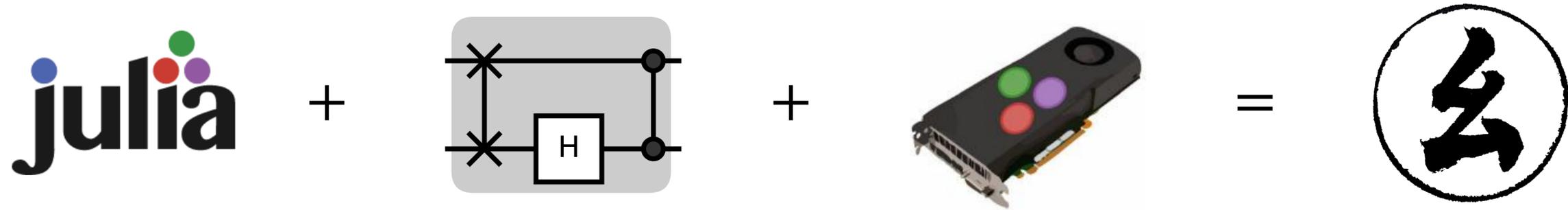
```

*Train a 10,000 layer,
300,000 parameter
circuit on a laptop*



<https://yaoquantum.org/>

Yao.jl: Extensible, Efficient Framework for Quantum Algorithm Design



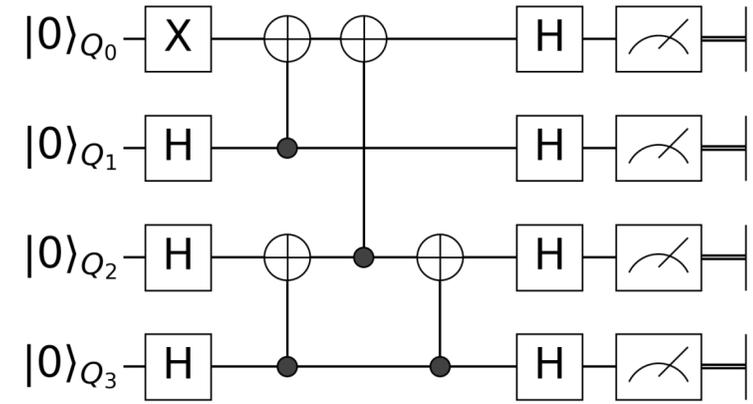
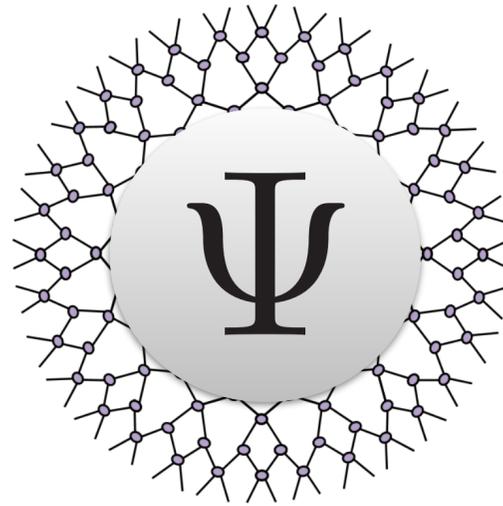
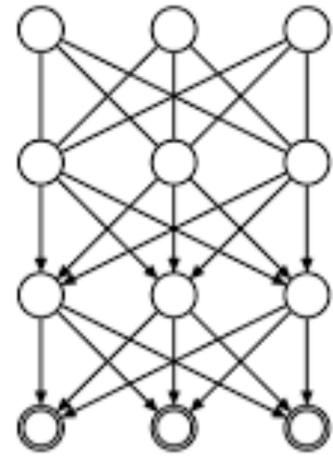
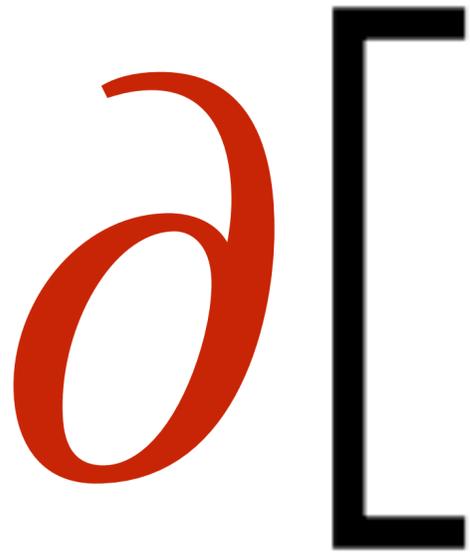
Xiu-Zhe Luo (IOP, CAS → Waterloo & PI)

Jin-Guo Liu (IOP, CAS → QuEra Computing & Harvard)

Features:

- Reversible AD engine for quantum circuits
- Batch parallelization with GPU acceleration
- Quantum block intermediate representation

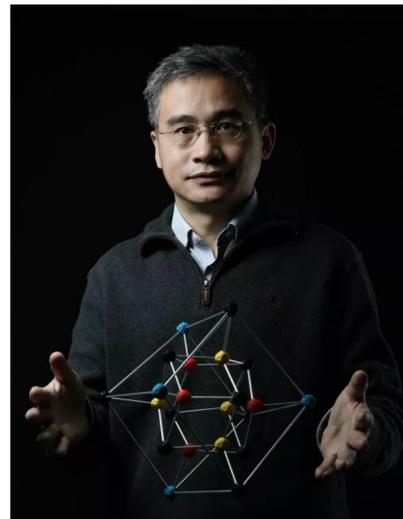




Thank you!



Hai-Jun Liao
IOP CAS



Tao Xiang
IOP CAS



Pan Zhang
ITP CAS



Jin-Guo Liu,
QuEra & Harvard



Xiu-Zhe Luo
Waterloo & PI