

Born-Infeld (BI) for AI:

Energy-Conserving Descent (ECD) for Optimization



G. Bruno De Luca*

Stanford Institute for Theoretical
Physics (SITP)



Eva Silverstein*

Acknowledgements \supset

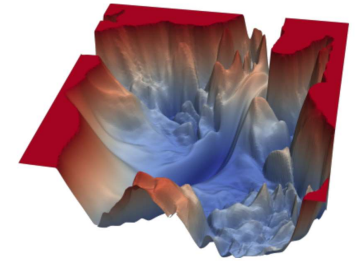
Discussions w/ J. Batson, Y. Kahn, D. Roberts on inflationary cosmology and optimization; early collaboration G. Panagopoulos, T. Bachlechner

New discussions/collaborations: ML (Kunin), Computational Chemistry (Zhang), Sampling (Robnik/Seljak), ...

Reviewers & conference organizers

Overview: Optimization of an objective function F

- Data analysis/Machine Learning [$F = \text{loss}$]
- Solving (Partial) Differential Equations
 $[F = \Sigma (\text{PDEs})^2 + (\text{boundary conditions})^2]$
- Many scientific applications



[Image from Li et al. , '18]

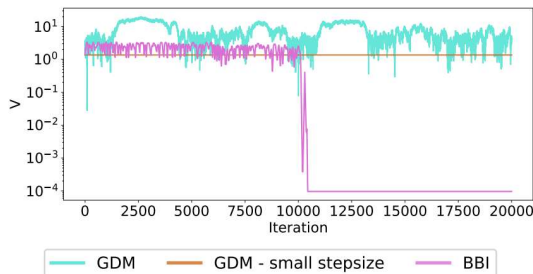
Gradient Descent with Momentum (GDM) can work well with modern tweaks.

Physical analogue: particle motion on potential energy $V = F$, *with friction*, discretized.

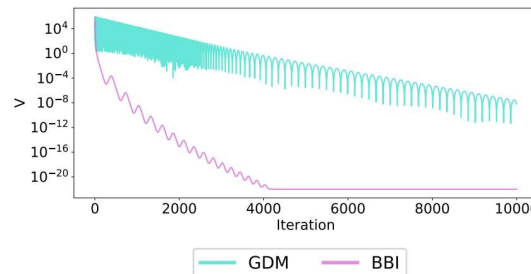
Our proposal: Energy Conserving Descent (**ECD**): discretized physical evolution, *without friction*, nonetheless slowing near minimal F . Examples include:

- BBI: relativistic, $(\text{speed limit})^2 = V = F - \Delta V$ [or more general $(\text{speed limit})^2 = g(V)$]
- Ruthless: non-relativistic, $\text{mass} \propto 1/g(V)$

Ackley 2d (nonconvex)



Zakharov 10d (shallow)



+ other synthetics, PDEs,
small ML (Cifar, MNIST,
Tiny ImageNet [new]),
chemistry, sampling [new]

No friction \Rightarrow **Energy Conservation** \Rightarrow favorable properties and improved calculability:
concrete formula for distribution of results: in all dims weighted toward small $V = F - \Delta V$

Physics of GDM

Particle descending a potential energy landscape V

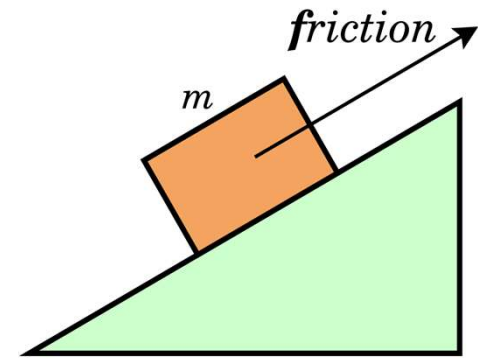
$$V(\Theta) = F(\Theta) - \Delta V$$

Familiar law of motion: Force = mass \times acceleration

$$-\nabla V - \mathbf{f}\dot{\Theta} = m\ddot{\Theta}$$

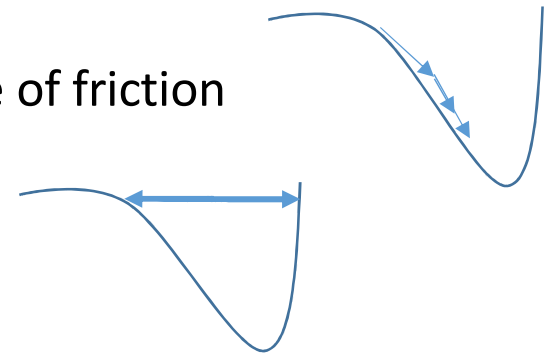
Friction coefficient $\mathbf{f} \Rightarrow$ Energy not conserved

First-order form: $p = m\dot{\Theta}$ $\dot{p} = -p\frac{\mathbf{f}}{m} - \nabla V$



Discretization \rightarrow GD with Momentum (GDM) + minibatches \rightarrow SGDM

- Energy $E = \frac{p^2}{2m} + V(\Theta)$ not conserved because of friction
- $\mathbf{f} = 0$ would conserve energy, but the particle flies quickly past $V \simeq 0$, spending very little time there (especially in high dimensions)



ECD: physical dynamics can conserve energy yet slow near $V=0$

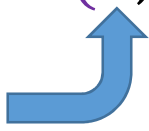
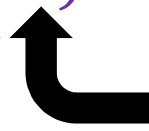
Next: explicit realizations

Explicit realizations of ECD

Change the dynamics to conserve Energy E and favor $V \simeq 0$

General

$$H(\Theta, \Pi) = E$$

Position vector (parameters)  Momentum vector 

Dynamical equations (cf. Newton's laws of motion): $\dot{\Theta} = \frac{\partial H}{\partial \Pi}$, $\dot{\Pi} = -\frac{\partial H}{\partial \Theta}$

1. **BI:** (speed limit)² = $V = F - \Delta V$, [or general function $g(V)$]

$$H = \sqrt{g(V)(\Pi^2 + g(V))} = g(V) / \sqrt{1 - \frac{\dot{\Theta}^2}{g(V)}}$$

- Cannot exceed relativistic speed limit: $\dot{\Theta}^2 \leq g(V)$ [ES, Tong, Alishahiha '04, cf. França et al. '20]

2. **Rootless (Ruthless):** mass $\propto 1/g(V)$ $H = \left(\frac{\Pi^2}{2m(V)}\right) = g(V) \Pi^2 = \frac{1}{2} m(V) \dot{\Theta}^2$

- Slows as the particle gets heavy: $m(V) \rightarrow \infty$, $g(V) \rightarrow 0 \Rightarrow \dot{\Theta}^2 \rightarrow 0$

Building ECD optimization algorithms

0. Choose the continuum dynamical system
1. Discretize the continuum equations of motion
 - e.g. BI with $g(V) = V$:
$$\sqrt{V(V + \vec{\pi}^2)} \equiv E$$
$$\pi_i(t + \Delta t) - \pi_i(t) = -\Delta t \frac{\partial_i V(\Theta(t))}{2} \left(\frac{E}{V} + \frac{V}{E} \right)$$
$$\theta_i(t + \Delta t) - \theta_i(t) = \Delta t \pi_i(t + \Delta t) \frac{V(\Theta(t))}{E}$$
2. Choose an initialization
 - Common choice: $\Pi(\theta) \Rightarrow E = V(\theta)$
 - Option: $E > \theta \Rightarrow$ choice of $\Pi(\theta)$ compatible with Energy eq.
3. Use discretized equation as update rules
4. Add other features
 - Enforce strict Energy conservation rescaling Π
 - Adaptive tuning of shift $\Delta V = F - V$ (next page)
 - Option: random rotation of momenta ("bouncing", explained later)
5. Test it!

DATA SET	SGD	BBI
MNIST	99.166 , 98.160	99.177 , 99.190
CIFAR-10	92.628 , 92.655	92.434 , 92.435

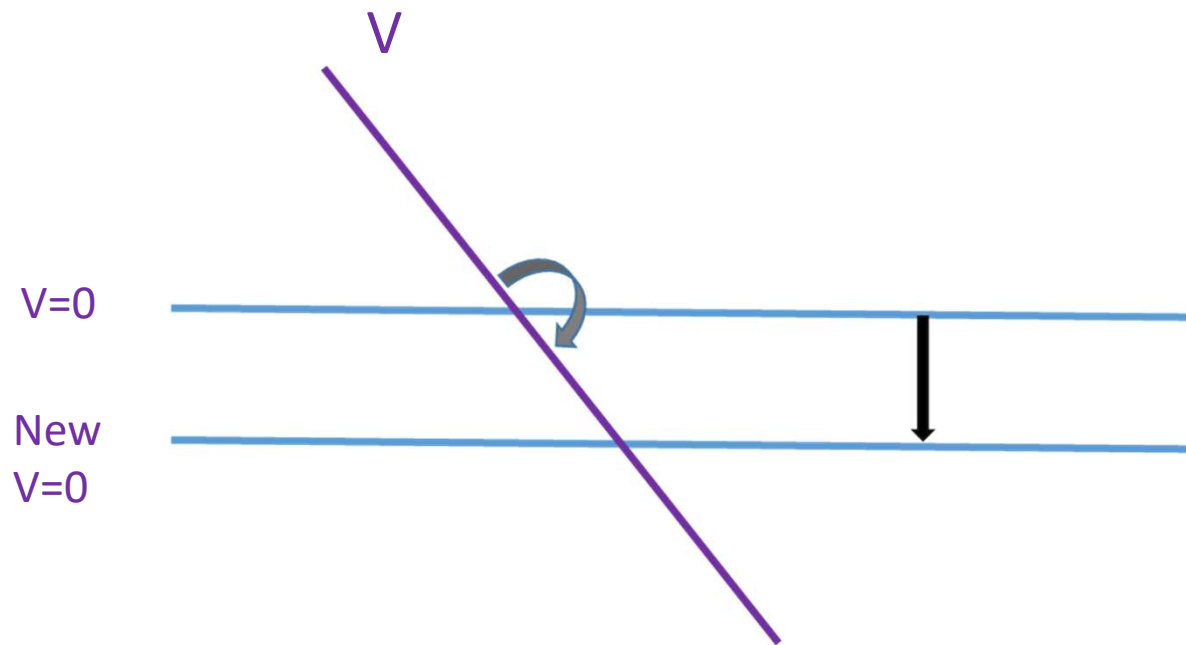
Modest (~50) statistics and limited hyper-parameter tuning (without all the tweaks on either side); just a check of basic competence. "Bouncing" not required here.

Automatic (adaptive) Tuning of ΔV

The value of the loss function F at the objective is not always known:

$$V = F - \Delta V$$

ΔV is a hyperparameter that can automatically adjust (recover from an over-estimate). **New** upgrade to optimizer code.



Given a too-high initial guess for ΔV , the loss extends to $V = F - \Delta V < 0$ and the trajectory will jump to a small negative value $V < 0$ due to the discreteness. Conditioned on this, ΔV may be lowered, iteratively tuning it.

Recap so far:

- Optimization of an objective function F
- Descent dynamics as (discrete) physical evolution on a potential $V = F - \Delta V$
- Equations of motion (update rules) obtained from a Hamiltonian H
 - Gradient Descent with Momentum: a time-dependent $H(\Pi, \Theta, t)$
 - Energy not conserved: $\dot{E} = -f \frac{\Pi^2}{m^2} \leq 0$
 - Simply removing friction ($f = 0$) does not converge
 - Alternative physics: Energy Conserving dynamical systems converging to $V \rightarrow 0$
$$E = H_{\text{ECD}}(\Theta, \Pi)$$
 - Energy is conserved: $\dot{E} = 0$
 - 2 explicit examples: **BI** [relativistic], **Ruthless** [$m = 1/g(V)$].
- Discretization gives update rules \rightarrow new optimization algorithms

Simple benchmarks show that the idea works: friction not needed for optimization.

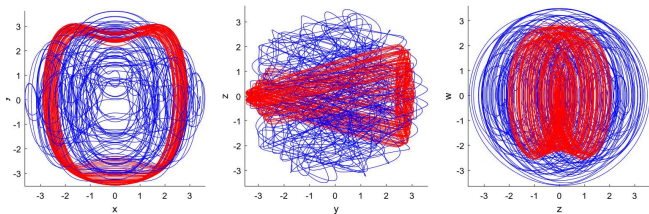
Next: advantages of conserving energy

Energy Conservation

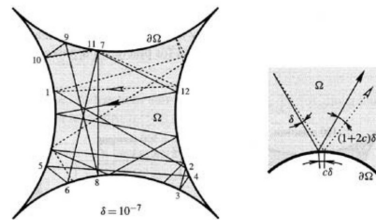
$$H = \frac{V}{\sqrt{1 - \frac{\dot{\theta}^2}{V}}} = \sqrt{V(V + \vec{\pi}^2)} \equiv E = \text{constant} \implies$$

- **Cannot** stop unless $V=E$ or $V=0$, so cannot stop in high local minimum

Can get stuck in orbit at high V . Generically such orbits are unstable: chaos – sensitive dependence on initial conditions – is typical in physical systems. Nearby trajectories disperse roughly on a *mixing* timescale.



[Image from Dong, Yuan, Du et al. '19]



[Image from Encyclopedia of Nonlinear Science, '04]

Chaos and *mixing* has been **proven** in mathematical billiards problems.

This inspires optional **Bounces** in BI algorithm above to reduce the *mixing time* \implies **BBI**

- *Phase space* (positions & momenta) *volume* is preserved under the evolution.
$$\text{Vol}(\text{phase space}) = \int d^n \Theta d^n \Pi \delta(H(\Pi, \Theta) - E)$$
- Past the *mixing time*, the probability to find a particle from a droplet (bundle of trajectories) in a region M of phase space is $\propto \text{Vol}(M)$

★ For ECD, phase space volume is strongly dominated near $V=0$:

$$\text{Vol}(\mathcal{M}) = \frac{2\pi^{n/2}}{\Gamma(n/2)} \int d^n\theta \int d\tilde{\pi} \tilde{\pi}^{n-1} \delta(\sqrt{V(V + \tilde{\pi}^2)} - E) = \frac{2\pi^{n/2}}{\Gamma(n/2)} \int d^n\theta \frac{E}{V} \left(\frac{E^2}{V} - V \right)^{\frac{n-2}{2}}$$

For $V \rightarrow 0$, $\text{Vol} \propto \int \frac{d^n\Theta}{V^{n/2}} = \int d\Omega \int d|\Theta| |\Theta|^{n-1} \frac{1}{V^{n/2}}$

For a basin $V \sim |\Theta|^2$, this becomes $\sim \int d\Omega \int d|\Theta| / |\Theta|$

$V \rightarrow g(V) \sim V^\eta \quad \eta > 1$ enhances the preference for $V=0$ (beats the effect of high dimension n !) ($g(V)$ also useful for sampling, in addition to optimization)

[GBDL, Roblik, Seljak, ES in progress]

- In contrast, pure momentum would not favor small V :

$$\text{Vol}(\mathcal{M}) \propto \int d^n\theta (E - V)^{\frac{n-2}{2}} \quad \text{frictionless non-relativistic momentum}$$

- The volume formula would not apply at all with friction (less predictive in that sense).

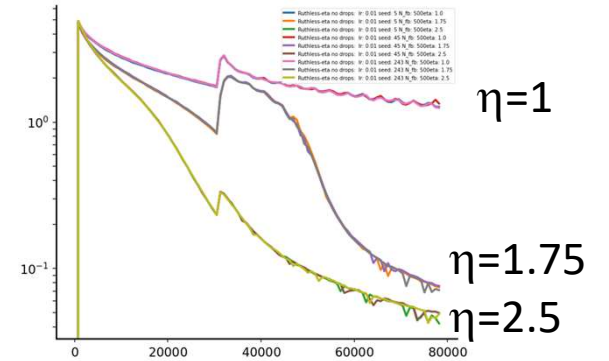
Exploiting the volume formula for image classification (preliminary)

- Enhancement of volume density for $\eta > 1$ near a quadratic minimum $V \sim \theta^2$:

$$\text{vol} \propto |\Theta|^{n(1-\eta)-1} d |\Theta|$$

- Small Tests on **Tiny-ImageNet*** with **D. Kunin** (+ImageNet 1K in progress)

Protocol: lr = 0.01, **no** lr drop needed, 500 bounces, Averaging of late-epoch weights (SWA) [Izmailov et al. '19]



Training loss decreases monotonically with η , improving test accuracy for intermediate $\eta > 1$

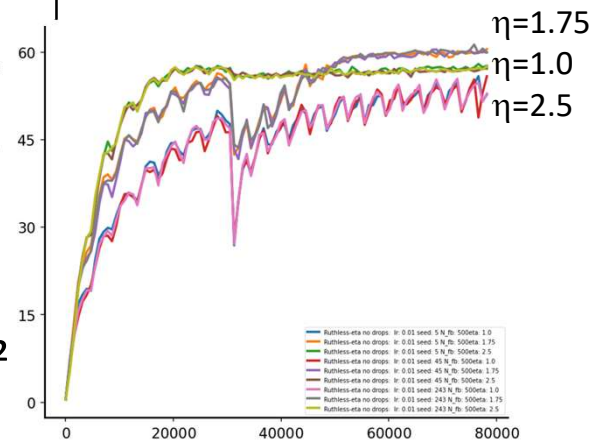
$m=1/V^\eta$	Accuracy	Accuracy (weights averaged)
$\eta=1$	55.44	62.12
$\eta=1.75$	61.3	64.1

Compared with SGD: with lr drops

(start 0.1, drop factor 0.1@ep. [30,60,80]) :

Accuracy: 62.52, Accuracy (weights averaged): 62.93

SGD: without lr drops is worse, as well as with loss \rightarrow loss²



[ECD also > best comparable SGDM in cf. Li et al. '21, Tanaka, Kunin et al. '20...]

*ResNet-18, epochs: 100, batch size: 128, weight decay: 10^{-4} , loss: Cross Entropy

Testing the volume formula

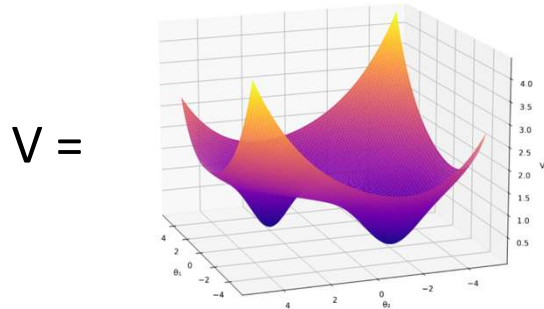
Evaluated in different regions predicts distribution of results (given mixing)

For $g(V) = V$:
$$Vol(\mathcal{M}_{\mathcal{I}}) = \frac{2\pi^{n/2}}{\Gamma(n/2)} E^{n-1} \int d^n(\theta - \theta_I) V^{-n/2}$$

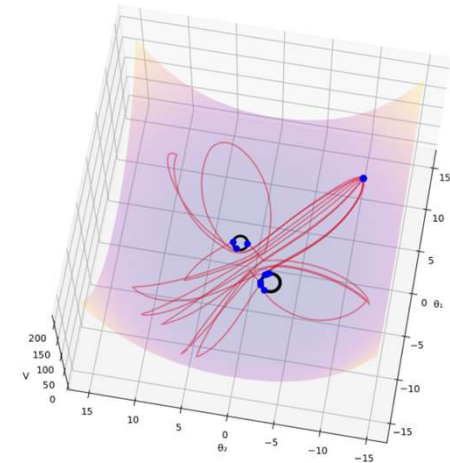
Near a minimum:

$$V \simeq V_I + \frac{1}{2} \sum_{i=1}^n m_{Ii}^2 (\theta_i - \theta_{Ii})^2 \quad Vol(\mathcal{M}_{\mathcal{I}}) \rightarrow b_n \left(\frac{2\pi^{n/2}}{\Gamma(n/2)} \right)^2 \frac{E^{n-1}}{\prod_i m_{Ii}} \log(V_I) \quad V_I \rightarrow 0$$

Empirical check:



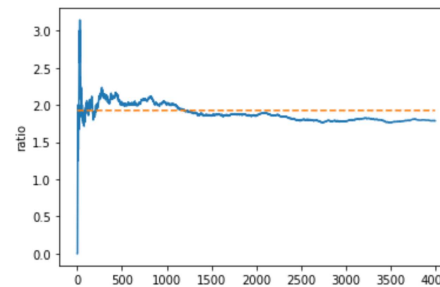
Bouncing trajectories find the 2 basins:



Prediction of ratio of convergence:

$$\frac{Vol(\mathcal{M}_1)}{Vol(\mathcal{M}_2)} \sim 1.93$$

Results:



Agreement within 10%

Figure 4: Partial ratios.

Behavior in shallow regions

Volume formula prefers **flatter** minima

ML lore: flatter minima generalize better

$$V \simeq V_I + \frac{1}{2} \sum_{i=1}^n m_{Ii}^2 (\theta_i - \theta_{Ii})^2$$

$$\text{Vol}(\mathcal{M}_I) \rightarrow b_n \left(\frac{2\pi^{n/2}}{\Gamma(n/2)} \right)^2 \frac{E^{n-1}}{\prod_i m_{Ii}} \log(V_I) \quad V_I \rightarrow 0, \quad m_{iI}^2 \rightarrow 0$$

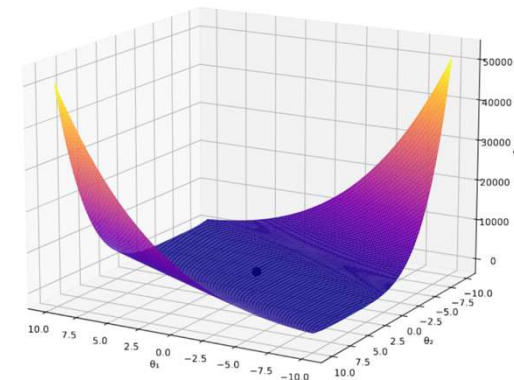
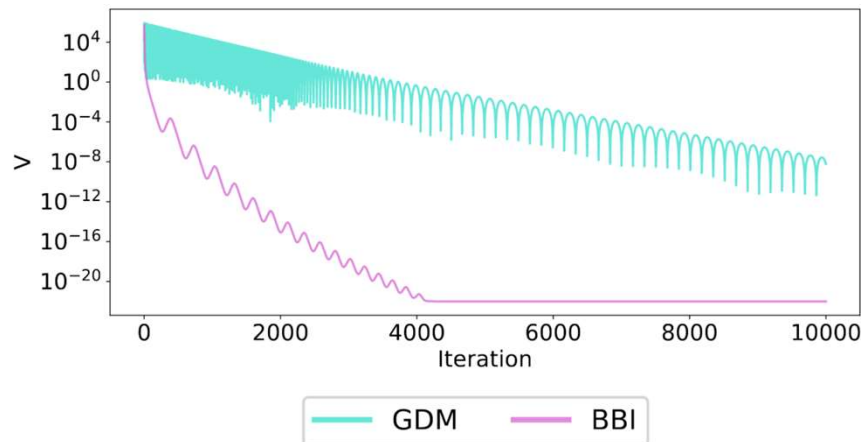
Prediction: BI is faster on **shallow** directions than GD

$$\Theta \sim e^{-mt/\sqrt{2}} \quad \text{vs} \quad \Theta \sim e^{-m^2 t/f}$$

Empirical check:

V = 10-dimensional Zakharov function

Results:



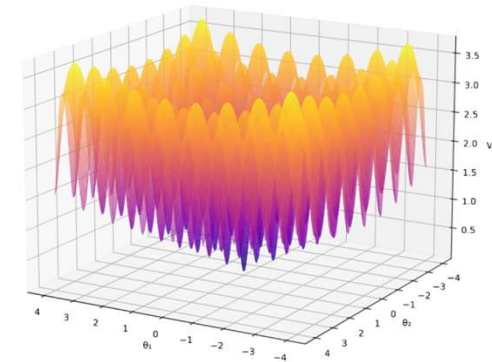
Hyperparameters tuned with hyperopt

Avoiding high local minima

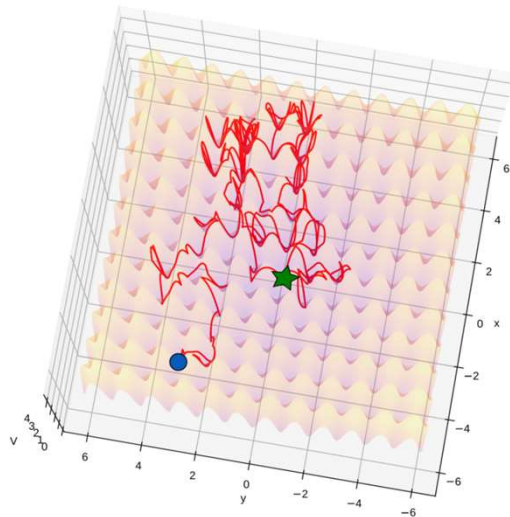
Energy conservation: **ECD** cannot stop in high local minima

Empirical check: Highly non-convex function

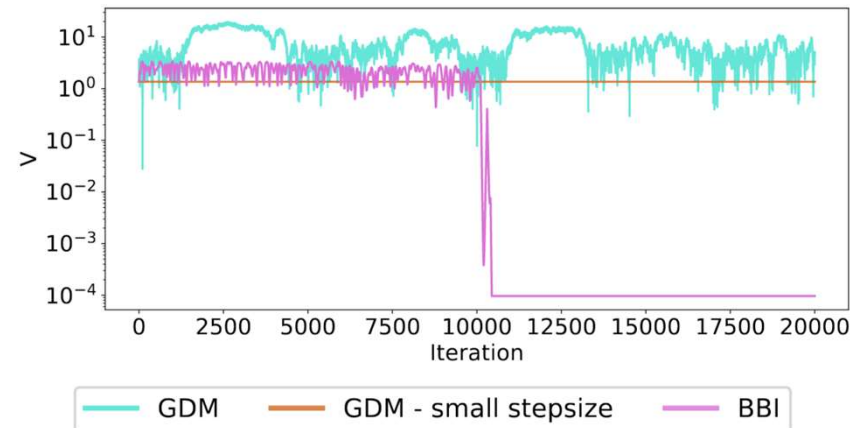
$V = 2$ -dim Ackley function :



Results:



BBI explores and finds the global minimum



Hyperoptimized fixed l_r , and for GDM also momentum. GDM either stuck in initial basin or helped out by 'catapult' mechanism [Lewkowycz et al. '20], , then more erratic (not settling in global minimum).

Summary comparison

ECD	FRICTION ((S)GDM, ...)
CONSERVES ENERGY E	FRICTION DRAINS E
CANNOT GET STUCK IN HIGH LOCAL MINIMUM	CAN STOP IN HIGH LOCAL MINIMUM
CANNOT OVERSHOOT $V = 0 = \nabla V$	CAN OVERSHOOT $V = 0 = \nabla V$
DEPENDS ON V AND ∇V	DEPENDS ONLY ON ∇V
ON SHALLOW REGION: $\theta \sim e^{-mt/\sqrt{2}}$	ON SHALLOW REGION: $\theta \sim e^{-m^2t/f}$
ANALYTIC PREDICTION FOR DISTRIBUTION	STOCHASTIC INTUITION FOR DISTRIBUTION
GENERALIZES	GENERALIZES

Generalization ok:
speed limit kicks in for
 $V \ll E$, Vol(phase space)
favors flat basins.

Statements persist with noise (mini-batches) in our prescription:
BBI speed limit tamps down noise, while the bounces (when needed) provide
controlled stochasticity for short mixing time.

Application: Solving Partial Differential Equations

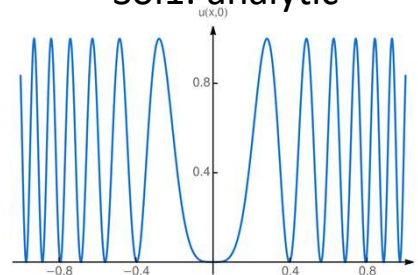
- Most common strategy with ML tools: a NN as ansatz for the PDE: [Lagaris et al. '98, ..., Raissi et al. '19,..]

$$F = V = \sum_{x \in \text{domain}} \text{PDE}[\mathcal{N}(x; \Theta)]^2 + \gamma \sum_{x \in \text{boundary}} \text{BC}[\mathcal{N}(x; \Theta)]^2 + R(\Theta)$$

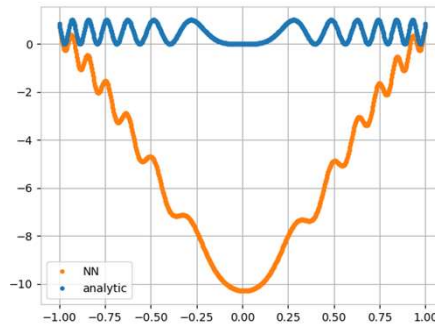
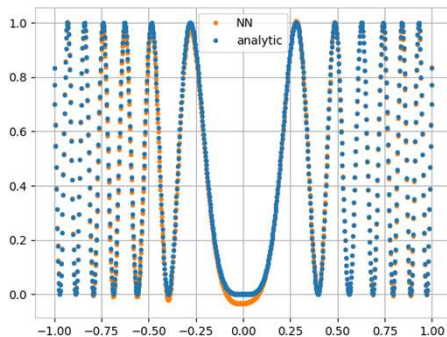
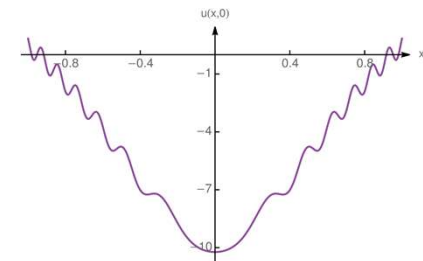
- We reverse-engineered hard (highly nonlinear) 2d PDEs with known multiple solution and checked if ECD optimization finds them

1d slices of *known* solutions:

Sol1: analytic



Sol2: numerical



1d slices of *learned* solutions

Found **both** from same initialization: bounces distribute results (mixing)

Ongoing work:

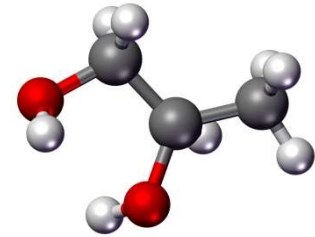
- Quantum Chemistry (with Zhang)
 - Find the minimum energy configuration of a molecule
 $F = \text{binding energy} < 0 \implies \text{requires } \Delta V$
 - Automatic tuning tested successfully
- Larger scale Machine Learning experiments (with Kunin)
 - Exploit the volume formula from frictionless dynamics for better generalization
- Efficient sampling from a function $\exp(-F)$ (with Robnik, Seljak)
 - Reverse engineer $g(V)$ such that

$$\text{Vol}(\text{phase space}) = \int d^n \Pi \int d^n \Theta \exp(-F) \delta(E - H(\Theta, \Pi)) \propto \int d^n \Theta \exp(-F)$$

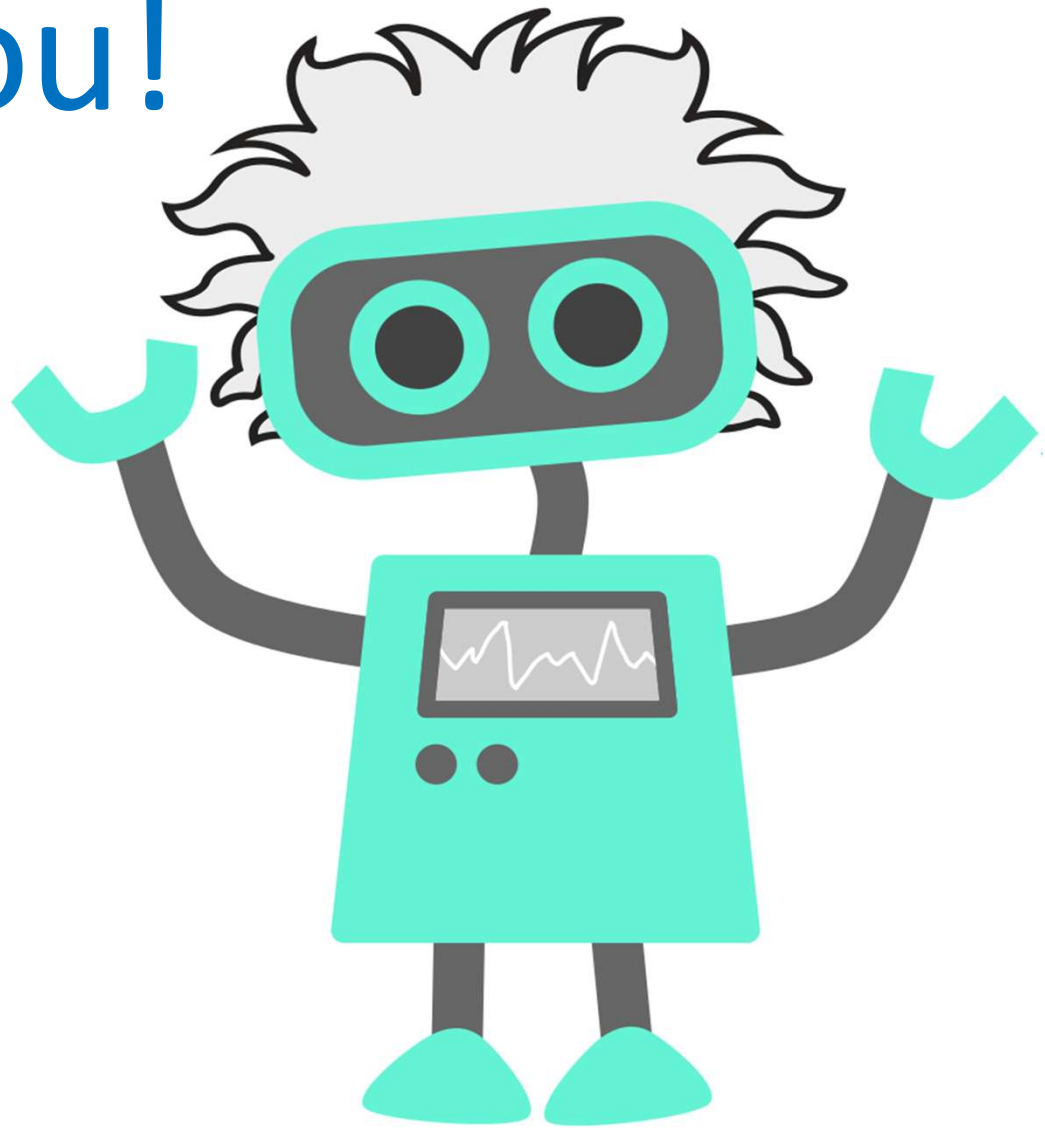
- In contrast to Hamiltonian Monte Carlo, no momentum sampling needed

Future directions:

- Feature learning theory and experiment
 - Bounces along the directions of hidden layer parameters



Thank you!



[Robot: publicdomainvectors.org,
Hair: He, Shuhan. (2020). Albert Einstein. Zenodo.
<https://doi.org/10.5281/zenodo.3926055>]