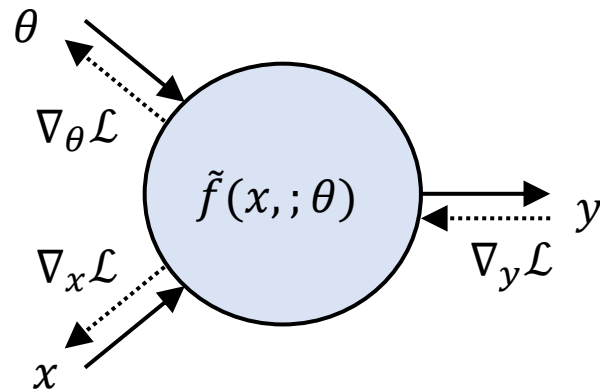# Exploiting Problem Structure in Deep Declarative Networks: Two Case Studies

Stephen Gould

Australian National University

The 1st International Workshop on Optimal Transport and Structured Data Modeling, AAAI, February 2022

# Deep declarative networks (DDNs)



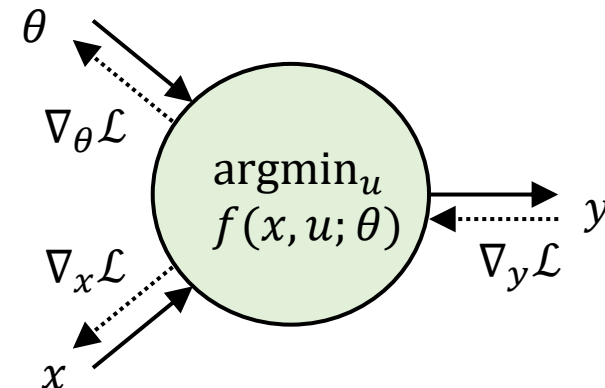In an **imperative node** the input-output relationship is explicitly defined as

$$y = \tilde{f}(x; \theta)$$

where $x$ is the input and $\theta$ are the parameters of the node.
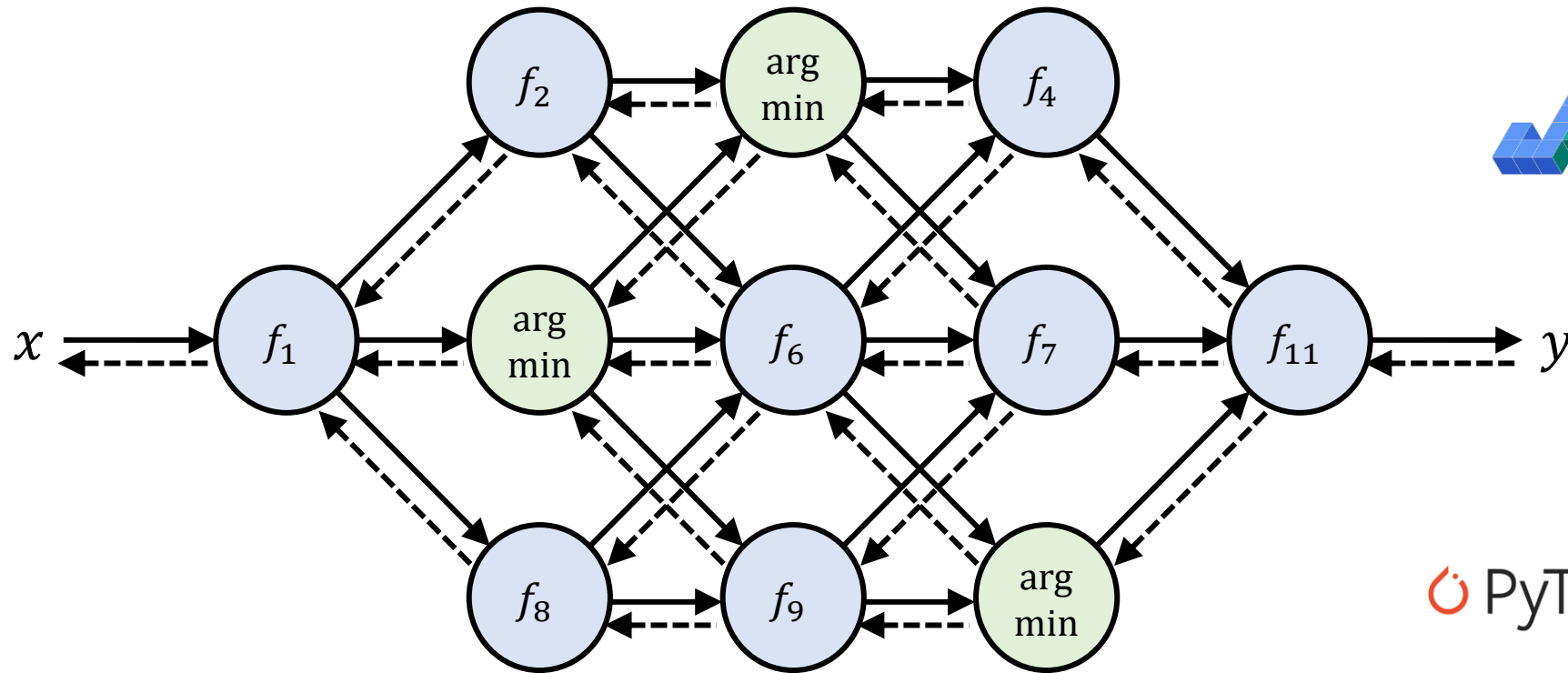
In a **declarative node** the input-output relationship is specified as the solution to an optimization problem

$$y \in \mathrm{argmin}_{u \in C} \, f(x, u; \theta)$$

where $f$ is the objective and $C$ are the constraints.

# Imperative and declarative nodes can co-exist

# Main technical question for DDNs

How do we compute $\dfrac{d}{dx} \text{argmin}_{u \in C(x)} f(x, u)$ ?

# Two answers

- Imperative approach: unroll the optimization procedure
  - **Advantages:** very simple, makes use of automatic differentiation so no additional coding is required
  - **Disadvantages:** need to store intermediate calculations in the forward pass, numerical issues when propagating through many iterations, potentially slow, may not be possible if non-differentiable steps are used in the forward pass

- Declarative approach: differentiate the optimality conditions to obtain closed-form expression for the gradient
  - Advantages/disadvantages to be continued…

# Main result for (smooth) DDNs

Given optimization problem parametrized by $x \in \mathbb{R}^n$, we define the solution set $Y(x) \colon \mathbb{R}^n \rightrightarrows \mathbb{R}^m$ as

smooth constraints

$$Y(x) = \text{argmin}_u \{ \underbrace{f(x, u)}_{\text{smooth objective}} : \overbrace{h(x, u) = 0_p, g(x, u) \leq 0_q}^{} \}$$

Then for any regular $y \in Y(x)$,

$$\frac{dy}{dx} = \underbrace{H^{-1} A^T (A H^{-1} A^T)^{-1} (A H^{-1} B - C) - H^{-1} B}_{\text{second partial derivatives of } f, h \text{ and } g}$$

# Simplified result

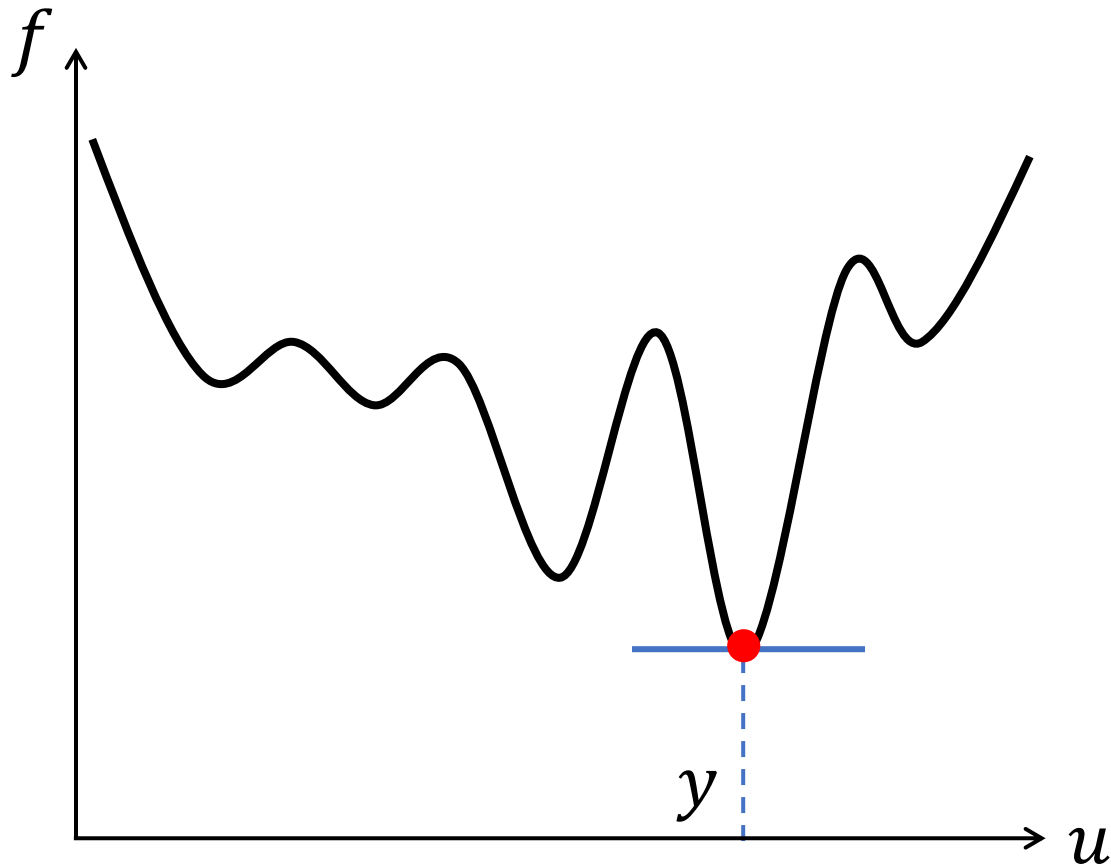Let $f: \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ be a twice differentiable function and let

$$y(x) \in \underset{u}{\mathrm{argmin}} \, \underbrace{f(x, u)}$$

<span style="color:blue">smooth objective</span>

The derivative of $f$ vanishes at $(x, y)$ so by Dini's implicit function theorem (1878)

$$\frac{dy(x)}{dx} = - \underbrace{\left(\frac{\partial^2 f}{\partial y^2}\right)^{-1} \frac{\partial^2 f}{\partial x \partial y}}$$

<span style="color:green">second partial derivatives of $f$</span>

# Proof sketch



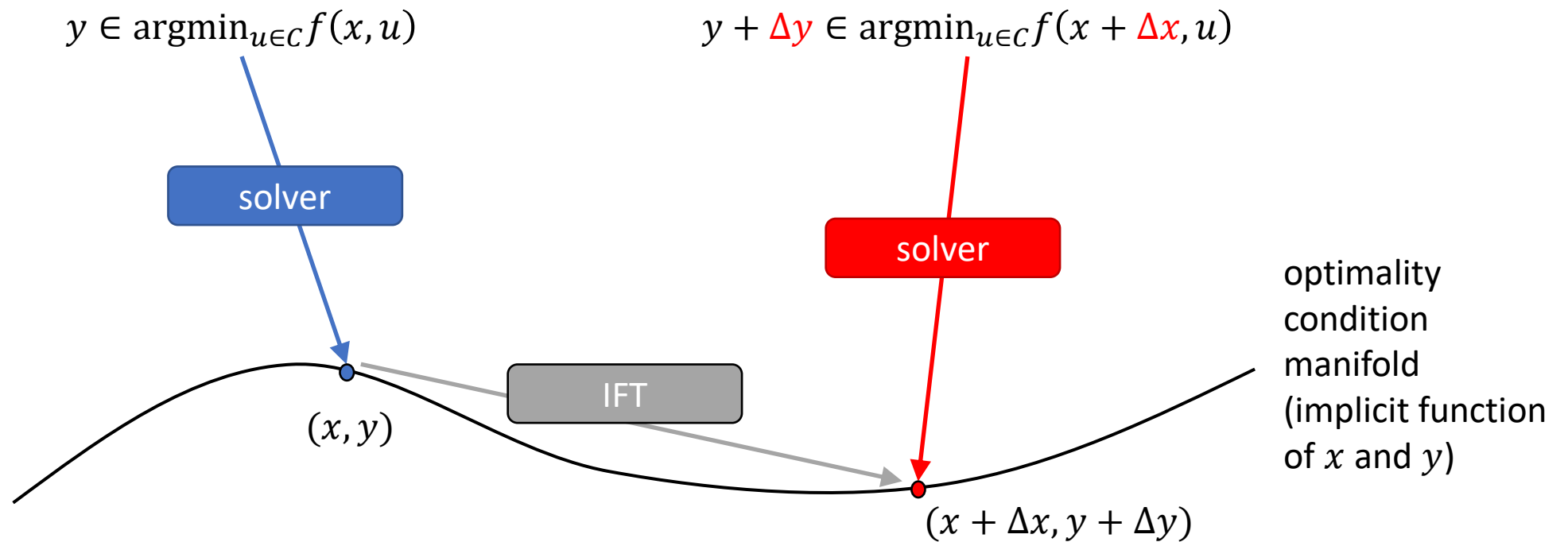$$y \in \operatorname{argmin}_u f(x, u) \Rightarrow \frac{\partial f(x, y)}{\partial y} = 0$$

LHS: $\quad \frac{d}{dx} \frac{\partial f(x,y)}{\partial y} = \frac{\partial^2 f(x,y)}{\partial x \partial y} + \frac{\partial^2 f(x,y)}{\partial y^2} \frac{dy}{dx}$

RHS: $\quad \frac{d}{dx} 0 = 0$

Rearranging gives $\dfrac{dy}{dx} = -\left(\dfrac{\partial^2 f}{\partial y^2}\right)^{-1} \dfrac{\partial^2 f}{\partial x \partial y}.$
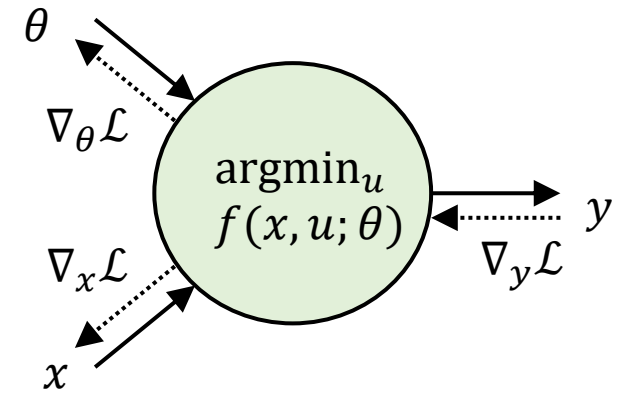
# Differentiable Optimization Concept

How does $y$ change as $x$ changes?

$$y \in \text{argmin}_{u \in C} f(x, u)$$

$$y + \Delta y \in \text{argmin}_{u \in C} f(x + \Delta x, u)$$

solver

solver

IFT

$(x, y)$

$(x + \Delta x, y + \Delta y)$

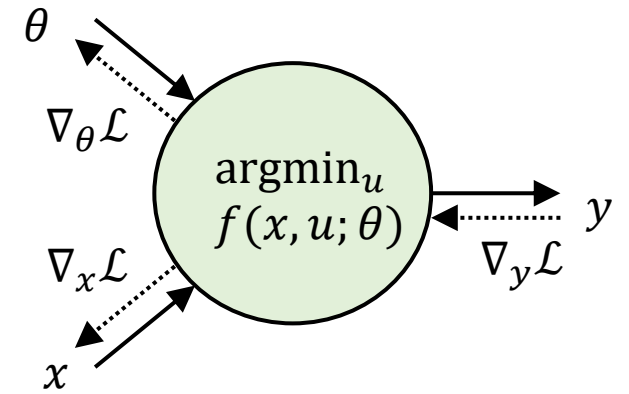optimality condition manifold (implicit function of $x$ and $y$)

# Backward pass summary

$$\frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}}{dy}\frac{dy}{dx}$$

# Backward pass summary

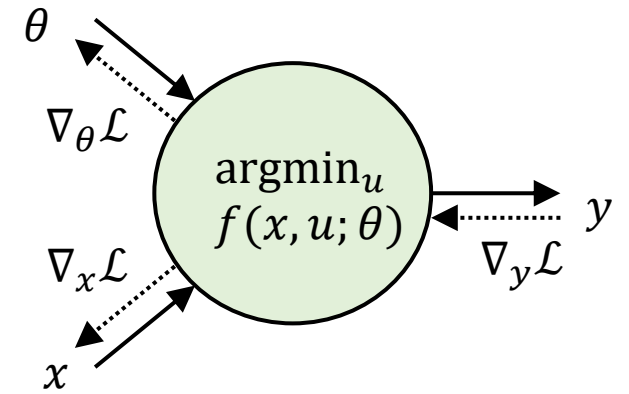$$\frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}}{dy}\frac{dy}{dx}$$



$$v^T\left(H^{-1}A^T\left(AH^{-1}A^T\right)^{-1}\left(AH^{-1}B - C\right) - H^{-1}B\right)$$

# Backward pass summary

$$\frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}}{dy}\frac{dy}{dx}$$



$$v^T\left(H^{-1}A^T\left(AH^{-1}A^T\right)^{-1}\left(AH^{-1}B - C\right) - H^{-1}B\right)$$

# Backward pass summary

$$\frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}}{dy}\frac{dy}{dx}$$



$$v^T \left( H^{-1}A^T \left( AH^{-1}A^T \right)^{-1} \left( AH^{-1}B - C \right) - H^{-1}B \right)$$

# Backward pass summary

$$\frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}}{dy}\frac{dy}{dx}$$



$$v^T\left(H^{-1}A^T\left(AH^{-1}A^T\right)^{-1}\left(AH^{-1}B - C\right) - H^{-1}B\right)$$

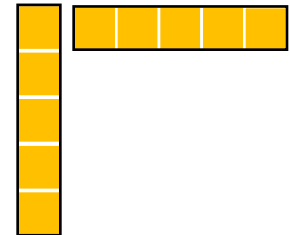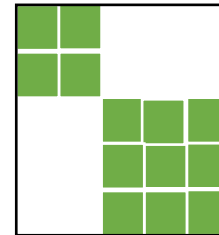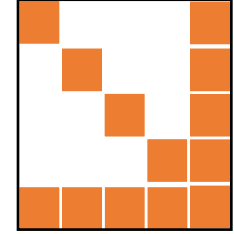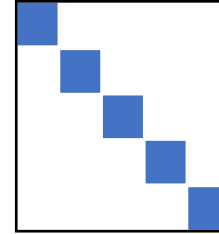$$\frac{d^2f}{dy^2} - \lambda\frac{d^2h}{dy^2} \qquad \frac{dh}{dy} \qquad \frac{dh}{dx} \qquad \frac{d^2f}{dydx} - \lambda\frac{d^2h}{dydx}$$

$$H^{-1}$$

# Exploiting structure

- Many problems exhibit structure that can be exploited for $H$ and other terms

- Our paper presents two case studies:
  - Robust vector pooling (omitted in this talk)
  - Optimal transport

# Entropic optimal transport

minimize    <span style="color:red">inner product</span> <span style="color:red">regularization</span>

$$\langle P, M \rangle + \frac{1}{\gamma} \mathrm{KL}(P \| rc^T)$$

subject to

$$P1 = r$$
$$P^T 1 = c$$

<span style="color:red">row sum and column sum constraints</span>

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \qquad \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \qquad \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

# Sinkhorn algorithm

**initialize** $P$ as $P_{ij} \leftarrow e^{-\gamma M_{ij}}$

**repeat** until convergence

    **for** $i$ in $1, \ldots, n$ **do**

        set $\alpha_i$ to the sum of the $i$-th row of $P$

        scale the $i$-th row of $P$ by $\dfrac{1}{\alpha_i}$

    **for** $i$ in $1, \ldots, n$ **do**

        set $\beta_j$ to the sum of the $j$-th column of $P$

        scale the $j$-th column of $P$ by $\dfrac{1}{\beta_j}$

**return** $P$

# Function mapping $M$ to $P$

$M$

m-by-n real
matrix

$P$

m-by-n positive
matrix

satisfying row
and column sum
constraints

# Backward pass for OT



$$\frac{d\mathcal{L}}{dM} = \frac{d\mathcal{L}}{dP}\overbrace{\frac{dP}{dM}}^{\text{m-by-n-by-m-by-n}}$$

$\underbrace{\phantom{\frac{d\mathcal{L}}{dM}}}_{\text{m-by-n}}$ $\underbrace{\phantom{\frac{dP}{dM}}}_{\text{m-by-n}}$

$$v^T(H^{-1}A^T(AH^{-1}A^T)^{-1}AH^{-1} - H^{-1})B$$
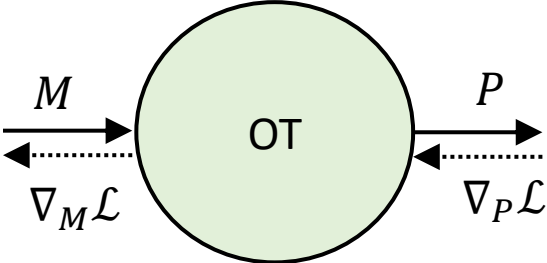
# Backward pass for OT



$$\underbrace{\frac{d\mathcal{L}}{dM}}_{\text{1-by-mn}} = \underbrace{\frac{d\mathcal{L}}{dP}}_{\text{1-by-mn}} \overbrace{\frac{dP}{dM}}^{\text{mn-by-mn}}$$

$$v^T \left( H^{-1} A^T (A H^{-1} A^T)^{-1} A H^{-1} - H^{-1} \right) B$$

# Computing $B$

$$f(P, M) = \sum_{i=1}^{m} \sum_{j=1}^{n} M_{ij} P_{ij} + \frac{1}{\gamma} \sum_{i=1}^{m} \sum_{j=1}^{n} P_{ij} \log \frac{P_{ij}}{r_i c_j}$$

$$B_{ij,kl} = \frac{\partial^2 f}{\partial P_{ij} \partial M_{kl}} = \begin{cases} 1 & \text{if } ij = kl \\ 0 & \text{otherwise} \end{cases}$$

# Computing $H^{-1}$

$$f(P,M) = \sum_{i=1}^{m}\sum_{j=1}^{n} M_{ij}P_{ij} + \frac{1}{\gamma}\sum_{i=1}^{m}\sum_{j=1}^{n} P_{ij}\log\frac{P_{ij}}{r_i c_j}$$

$$H^{-1} = \left(\left[\frac{\partial^2 f}{\partial P_{ij}\partial P_{kl}}\right]_{ij,kl}\right)^{-1} = \gamma\,\mathrm{diag}\big(\mathrm{vec}(P)\big)$$

# Computing $AH^{-1}A^T$

row and column sum constraints

$$A = \begin{bmatrix} 0_n^T & 1_n^T & \dots & 0_n^T \\ \vdots & \vdots & \ddots & \vdots \\ 0_n^T & 0_n^T & \dots & 1_n^T \\ I_{n \times n} & I_{n \times n} & \dots & I_{n \times n} \end{bmatrix}$$

$$AH^{-1}A^T = \gamma \begin{bmatrix} \text{diag}(r_{2:m}) & P_{2:m,1:n} \\ P_{2:m,1:n}^T & \text{diag}(c) \end{bmatrix}$$

# Inverting $AH^{-1}A^T$

$$\left(AH^{-1}A^T\right)^{-1} = \frac{1}{\gamma}\begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix}$$

$$\Lambda_{11} = \left(\text{diag}(r_{2:m}) - P_{2:m,1:n}\text{diag}(c)^{-1}P_{2:m,1:n}^T\right)^{-1}$$

$$\Lambda_{12} = -\Lambda_{11}P_{2:m,1:n}\text{diag}(c)^{-1}$$

$$\Lambda_{22} = \text{diag}(c)^{-1} - \text{diag}(c)^{-1}P_{2:m,1:n}^T\Lambda_{12}$$

$$\gamma v^T \vec{P} \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \vec{P} - \gamma v^T \vec{P}$$

```
# initialize backward gradients (-v^T H^{-1} B)
dJdM = -1.0 * gamma * P * dJdP

# compute [vHAt1, vHAt2] = -v^T H^{-1} A^T
vHAt1 = torch.sum(dJdM[:, 1:m, 0:n], dim=2)
vHAt2 = torch.sum(dJdM, dim=1)

# compute [v1, v2] = -v^T H^{-1} A^T (A H^{-1] A^T)^{-1}
P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
block_11 = torch.cholesky(torch.diag_embed(r[:, 1:m]) - torch.einsum("bij,bkj->bik", P[:, 1:m, 0:n], P_over_c))
block_12 = torch.cholesky_solve(P_over_c, block_11)
block_22 = torch.diag_embed(1.0 / c) + torch.einsum("bji,bjk->bik", block_12, P_over_c)

v1 = torch.cholesky_solve(vHAt1.view(batches, m-1, 1), block_11).view(batches, m-1) - torch.einsum("bi,bji->bj", vHAt2, block_12)
v2 = torch.einsum("bi,bij->bj", vHAt2, block_22) - torch.einsum("bi,bij->bj", vHAt1, block_12)

# compute v^T H^{-1} A^T (A H^{-1] A^T)^{-1} A H^{-1} B - v^T H^{-1} B
dJdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
dJdM -= v2.view(batches, 1, n) * P
```

$$\gamma v^T \vec{P} \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \vec{P} - \gamma v^T \vec{P}$$

```
# initialize backward gradients (-v^T H^{-1} B)
dJdM = -1.0 * gamma * P * dJdP
```

```
# compute [vHAt1, vHAt2] = -v^T H^{-1} A^T
vHAt1 = torch.sum(dJdM[:, 1:m, 0:n], dim=2)
vHAt2 = torch.sum(dJdM, dim=1)

# compute [v1, v2] = -v^T H^{-1} A^T (A H^{-1] A^T)^{-1}
P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
block_11 = torch.cholesky(torch.diag_embed(r[:, 1:m]) - torch.einsum("bij,bkj->bik", P[:, 1:m, 0:n], P_over_c))
block_12 = torch.cholesky_solve(P_over_c, block_11)
block_22 = torch.diag_embed(1.0 / c) + torch.einsum("bji,bjk->bik", block_12, P_over_c)

v1 = torch.cholesky_solve(vHAt1.view(batches, m-1, 1), block_11).view(batches, m-1) - torch.einsum("bi,bji->bj", vHAt2, block_12)
v2 = torch.einsum("bi,bij->bj", vHAt2, block_22) - torch.einsum("bi,bij->bj", vHAt1, block_12)

# compute v^T H^{-1} A^T (A H^{-1] A^T)^{-1} A H^{-1} B - v^T H^{-1} B
dJdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
dJdM -= v2.view(batches, 1, n) * P
```

$$\gamma v^T \vec{P} \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \vec{P} - \gamma v^T \vec{P}$$

```
# initialize backward gradients (-v^T H^{-1} B)
dJdM = -1.0 * gamma * P * dJdP

# compute [vHAt1, vHAt2] = -v^T H^{-1} A^T
vHAt1 = torch.sum(dJdM[:, 1:m, 0:n], dim=2)
vHAt2 = torch.sum(dJdM, dim=1)

# compute [v1, v2] = -v^T H^{-1} A^T (A H^{-1] A^T)^{-1}
P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
block_11 = torch.cholesky(torch.diag_embed(r[:, 1:m]) - torch.einsum("bij,bkj->bik", P[:, 1:m, 0:n], P_over_c))
block_12 = torch.cholesky_solve(P_over_c, block_11)
block_22 = torch.diag_embed(1.0 / c) + torch.einsum("bji,bjk->bik", block_12, P_over_c)

v1 = torch.cholesky_solve(vHAt1.view(batches, m-1, 1), block_11).view(batches, m-1) - torch.einsum("bi,bji->bj", vHAt2, block_12)
v2 = torch.einsum("bi,bij->bj", vHAt2, block_22) - torch.einsum("bi,bij->bj", vHAt1, block_12)

# compute v^T H^{-1} A^T (A H^{-1] A^T)^{-1} A H^{-1} B - v^T H^{-1} B
dJdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
dJdM -= v2.view(batches, 1, n) * P
```

$$\gamma v^T \vec{P} \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \vec{P} - \gamma v^T \vec{P}$$

```
# initialize backward gradients (-v^T H^{-1} B)
dJdM = -1.0 * gamma * P * dJdP

# compute [vHAt1, vHAt2] = -v^T H^{-1} A^T
vHAt1 = torch.sum(dJdM[:, 1:m, 0:n], dim=2)
vHAt2 = torch.sum(dJdM, dim=1)

# compute [v1, v2] = -v^T H^{-1} A^T (A H^{-1] A^T)^{-1}
P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
block_11 = torch.cholesky(torch.diag_embed(r[:, 1:m]) - torch.einsum("bij,bkj->bik", P[:, 1:m, 0:n], P_over_c))
block_12 = torch.cholesky_solve(P_over_c, block_11)
block_22 = torch.diag_embed(1.0 / c) + torch.einsum("bji,bjk->bik", block_12, P_over_c)

v1 = torch.cholesky_solve(vHAt1.view(batches, m-1, 1), block_11).view(batches, m-1) - torch.einsum("bi,bji->bj", vHAt2, block_12)
v2 = torch.einsum("bi,bij->bj", vHAt2, block_22) - torch.einsum("bi,bij->bj", vHAt1, block_12)

# compute v^T H^{-1} A^T (A H^{-1] A^T)^{-1} A H^{-1} B - v^T H^{-1} B
dJdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
dJdM -= v2.view(batches, 1, n) * P
```

$$\Lambda_{11} = \left(\mathrm{diag}(r_{2:m}) - P_{2:m,1:n}\mathrm{diag}(c)^{-1}P_{2:m,1:n}^T\right)^{-1}$$
$$= (LL^T)^{-1}$$

block_11 = torch.cholesky(…)

$$\Lambda_{12} = -\Lambda_{11}P_{2:m,1:n}\mathrm{diag}(c)^{-1}$$
$$= -L^{-T}L^{-1}P_{2:m,1:n}\mathrm{diag}(c)^{-1}$$

block_12 = torch.cholesky_solve(…, block_11)

$$\gamma v^T \vec{P} [A_1^T \quad A_2^T] \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \vec{P} - \gamma v^T \vec{P}$$

```
# initialize backward gradients (-v^T H^{-1} B)
dJdM = -1.0 * gamma * P * dJdP

# compute [vHAt1, vHAt2] = -v^T H^{-1} A^T
vHAt1 = torch.sum(dJdM[:, 1:m, 0:n], dim=2)
vHAt2 = torch.sum(dJdM, dim=1)

# compute [v1, v2] = -v^T H^{-1} A^T (A H^{-1] A^T)^{-1}
P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
block_11 = torch.cholesky(torch.diag_embed(r[:, 1:m]) - torch.einsum("bij,bkj->bik", P[:, 1:m, 0:n], P_over_c))
block_12 = torch.cholesky_solve(P_over_c, block_11)
block_22 = torch.diag_embed(1.0 / c) + torch.einsum("bji,bjk->bik", block_12, P_over_c)
```

```
v1 = torch.cholesky_solve(vHAt1.view(batches, m-1, 1), block_11).view(batches, m-1) - torch.einsum("bi,bji->bj", vHAt2, block_12)
v2 = torch.einsum("bi,bij->bj", vHAt2, block_22) - torch.einsum("bi,bij->bj", vHAt1, block_12)
```

```
# compute v^T H^{-1} A^T (A H^{-1} A^T)^{-1} A H^{-1} B - v^T H^{-1} B
dJdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
dJdM -= v2.view(batches, 1, n) * P
```

$$\gamma v^T \vec{P} \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \vec{P} - \gamma v^T \vec{P}$$

```
# initialize backward gradients (-v^T H^{-1} B)
dJdM = -1.0 * gamma * P * dJdP

# compute [vHAt1, vHAt2] = -v^T H^{-1} A^T
vHAt1 = torch.sum(dJdM[:, 1:m, 0:n], dim=2)
vHAt2 = torch.sum(dJdM, dim=1)

# compute [v1, v2] = -v^T H^{-1} A^T (A H^{-1] A^T)^{-1}
P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
block_11 = torch.cholesky(torch.diag_embed(r[:, 1:m]) - torch.einsum("bij,bkj->bik", P[:, 1:m, 0:n], P_over_c))
block_12 = torch.cholesky_solve(P_over_c, block_11)
block_22 = torch.diag_embed(1.0 / c) + torch.einsum("bji,bjk->bik", block_12, P_over_c)

v1 = torch.cholesky_solve(vHAt1.view(batches, m-1, 1), block_11).view(batches, m-1) - torch.einsum("bi,bji->bj", vHAt2, block_12)
v2 = torch.einsum("bi,bij->bj", vHAt2, block_22) - torch.einsum("bi,bij->bj", vHAt1, block_12)

# compute v^T H^{-1} A^T (A H^{-1] A^T)^{-1} A H^{-1} B - v^T H^{-1} B
dJdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
dJdM -= v2.view(batches, 1, n) * P
```
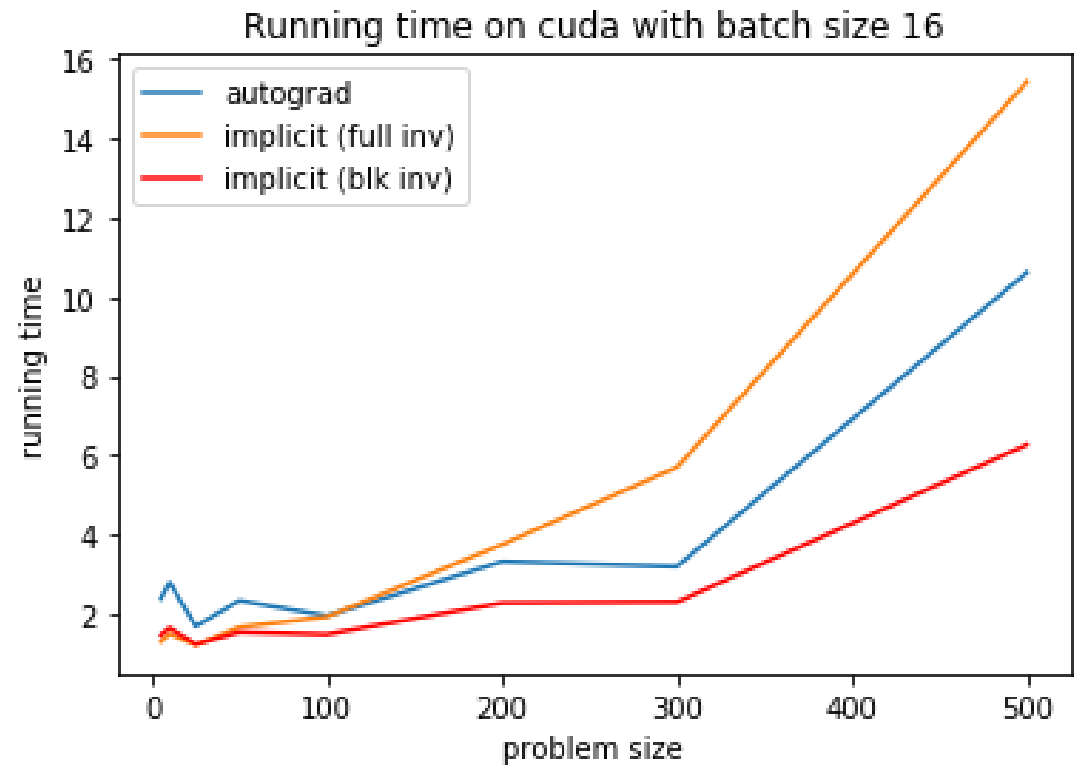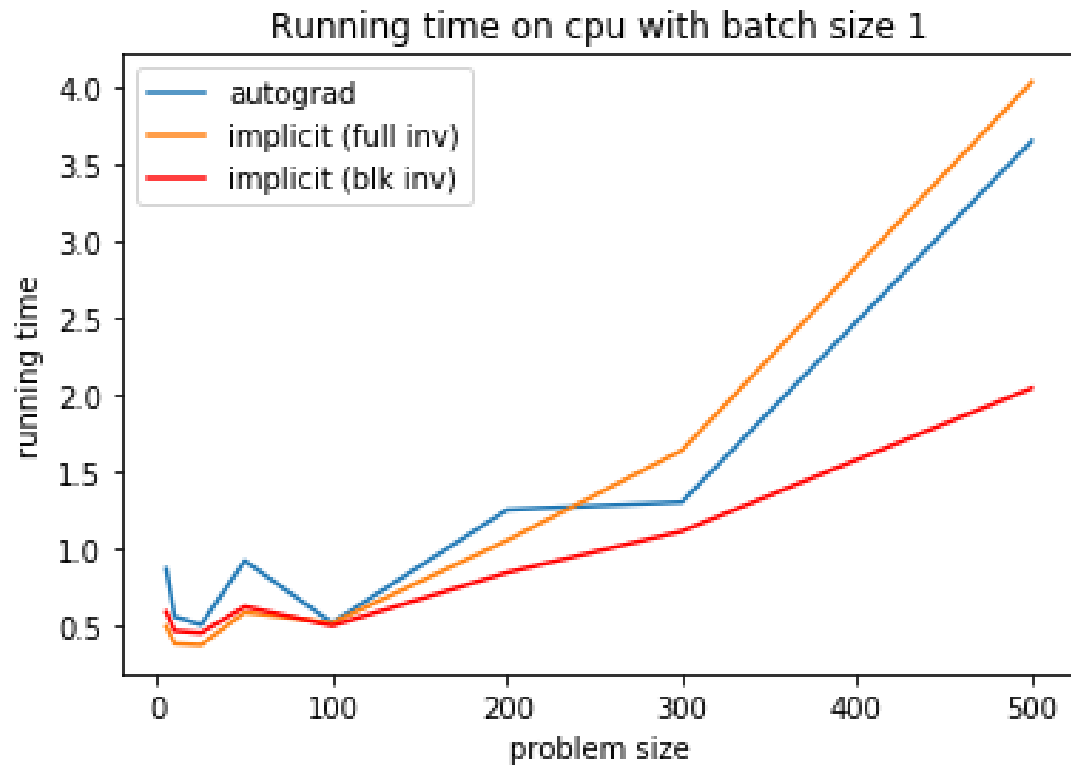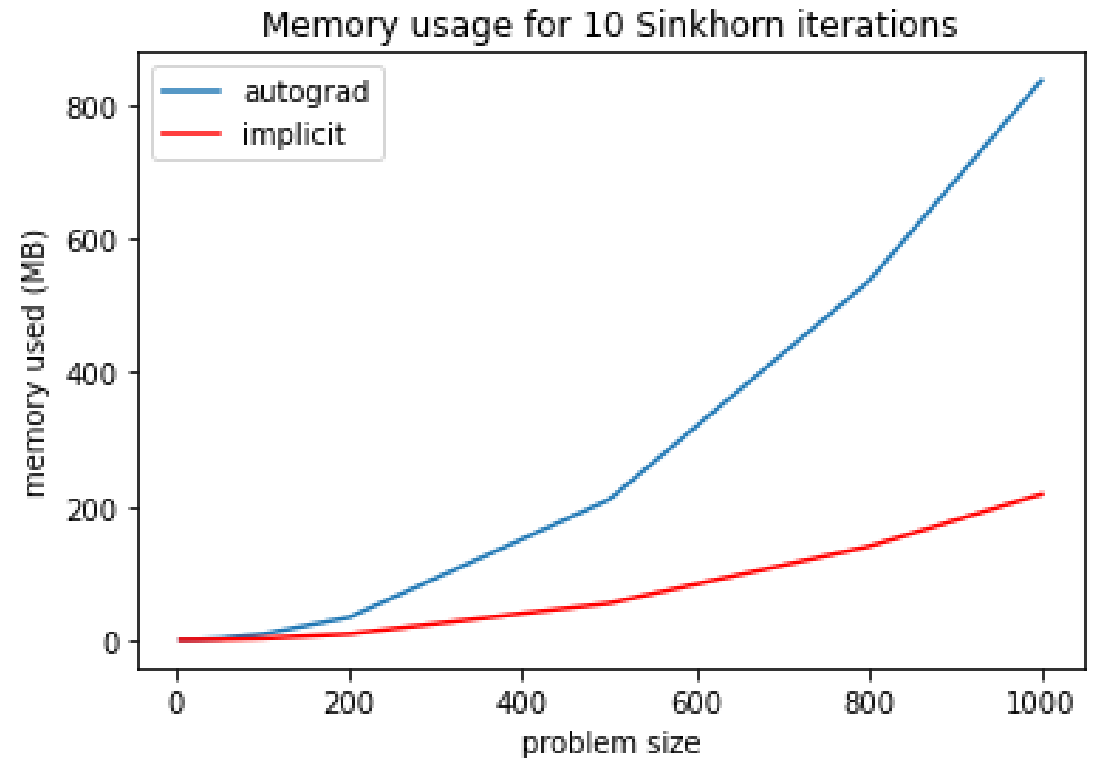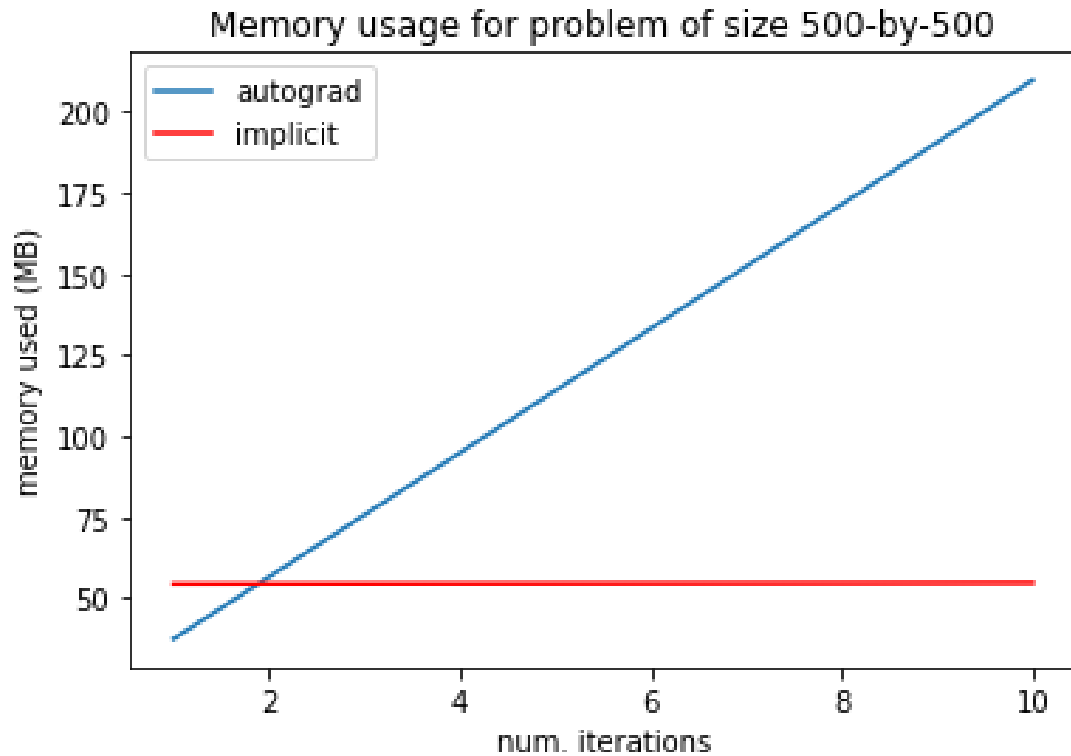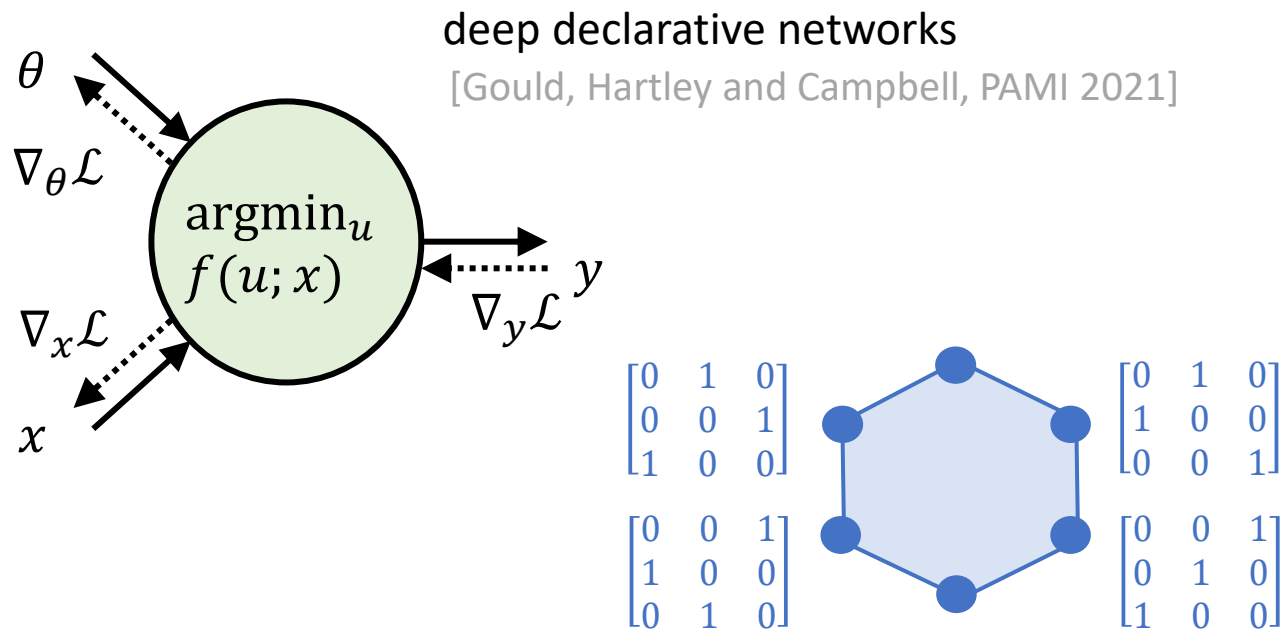
# Unrolling vs implicit differentiation

speed

# Unrolling vs implicit differentiation

memory

# Summary & Questions

deep declarative networks

[Gould, Hartley and Campbell, PAMI 2021]

$$\text{argmin}_u \, f(u; x)$$

$\theta$

$\nabla_\theta \mathcal{L}$

$\nabla_x \mathcal{L}$

$x$

$\nabla_y \mathcal{L}$ $y$

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

exploiting problem structure in DDNs: two case studies

[Gould, Campbell, Ben-Shabat, Koneputugodage and Xu, OT-SDM@AAAI 2022]

code and tutorials at
http://deepdeclarativenetworks.com

collaborators