Geometric data analysis, beyond convolutions

Jean Feydy King's College London – December 2019

Imperial College London Collaboration with B. Charlier, J. Glaunès (KeOps library); F.-X. Vialard, G. Peyré, T. Séjourné, A. Trouvé (OT theory).

• Typical French curriculum (ENS Maths + MVA + Siemens)

- Typical French curriculum (ENS Maths + MVA + Siemens)
- 2016-2019: PhD with Alain Trouvé (ENS Paris-Saclay)

- Typical French curriculum (ENS Maths + MVA + Siemens)
- 2016-2019: PhD with Alain Trouvé (ENS Paris-Saclay)
- 2019-2022: PostDoc at Imperial, Dept. of Computing

- Typical French curriculum (ENS Maths + MVA + Siemens)
- 2016-2019: PhD with Alain Trouvé (ENS Paris-Saclay)
- 2019-2022: PostDoc at Imperial, Dept. of Computing

- Typical French curriculum (ENS Maths + MVA + Siemens)
- 2016-2019: PhD with Alain Trouvé (ENS Paris-Saclay)
- 2019-2022: PostDoc at Imperial, Dept. of Computing

 \implies Looking for contacts in London!



Sensor data







Computational anatomy [CSG19, LSG⁺18, CMN14]

Three main problems:



Spot patterns

Analyze variations

Fit models

2010-2020: the deep learning revolution [Mal16, PMC11, WZTF17]

Explicit wavelet transform



2010-2020: the deep learning revolution [Mal16, PMC11, WZTF17]

Explicit wavelet transform \longrightarrow data-driven CNN



2010-2020: the deep learning revolution [Mal16, PMC11, WZTF17]

Explicit wavelet transform \longrightarrow data-driven CNN \longrightarrow deep CNN



(Wavelets \rightarrow CNNs) = improvement for... [NN16, Ola18]



Texture processing

(Wavelets \rightarrow CNNs) = improvement for... [NN16, Ola18]



Texture processing

Feature detection

Segmentation with U-nets [RFB15]



Architecture

Input

Output



 $\mathsf{Advection} \neq \mathsf{Convolution}$



Advection \neq Convolution



 $\mathsf{Advection} \neq \mathsf{Convolution}$



 $\mathsf{Advection} \neq \mathsf{Convolution}$



 $\mathsf{Advection} \neq \mathsf{Convolution}$



 $\mathsf{Advection} \neq \mathsf{Convolution}$





 $\mathsf{Advection} \neq \mathsf{Convolution}$





 $\mathsf{Advection} \neq \mathsf{Convolution}$





 $\mathsf{Advection} \neq \mathsf{Convolution}$

Mesh deformation

 \Longrightarrow Geometry is easier with **coordinates**: point clouds and meshes.



 $\mathsf{Advection} \neq \mathsf{Convolution}$

Mesh deformation

 \implies Geometry is easier with **coordinates**: point clouds and meshes. **Problem:** not supported well by TensorFlow and PyTorch.

- Efficient GPU implementations
 - \longrightarrow KeOps extension for PyTorch, NumPy, Matlab, R.

- Efficient GPU implementations
 - \longrightarrow KeOps extension for PyTorch, NumPy, Matlab, R.
- Geometric loss functions
 - \longrightarrow Optimal transport = Wasserstein distance.

- Efficient GPU implementations
 - \longrightarrow KeOps extension for PyTorch, NumPy, Matlab, R.
- Geometric loss functions
 - \longrightarrow Optimal transport = Wasserstein distance.
- Robust deformation architectures
 - \longrightarrow After the break :--)

Combining graphics and deep learning



Nvidia RTX 2080 Ti = 4,352 cores, 11Gb RAM for ~1,500\$.

Designed for graphics: $1\,{ m GF}{ m lop}\simeq 1\,{ m \$}$





Today

Theano (RIP), TensorFlow and PyTorch combine:

• Matlab-like Python interface.

Theano (RIP), TensorFlow and PyTorch combine:

- Matlab-like Python interface.
- CPU and GPU backends.
- Matlab-like Python interface.
- CPU and GPU backends.
- Automatic differentiation engine.

- Matlab-like Python interface.
- CPU and GPU backends.
- Automatic differentiation engine.
- Support for imaging (convolutions) and linear algebra.

- Matlab-like Python interface.
- CPU and GPU backends.
- Automatic differentiation engine.
- Support for imaging (convolutions) and linear algebra.

- Matlab-like Python interface.
- CPU and GPU backends.
- Automatic differentiation engine.
- Support for imaging (convolutions) and linear algebra.
- \implies Ideally suited for research!

PyTorch, in practice

```
# Store arbitrary arrays on the GPU:
N, D = 10000, 3
x = torch.rand(N, D) # Point clouds in [0,1]<sup>3</sup>
y = torch.rand(N, D) # encoded as N-by-D arrays
```

PyTorch, in practice

```
# Store arbitrary arrays on the GPU:
N, D = 10000, 3
x = torch.rand(N, D) # Point clouds in [0,1]<sup>3</sup>
y = torch.rand(N, D) # encoded as N-by-D arrays
# Re-indexing:
x_i = x[:,None,:] # shape (N,D) -> (N,1,D)
y_j = y[None,:,:] # shape (N,D) -> (1,N,D)
```

PyTorch, in practice

```
# Store arbitrary arrays on the GPU:
N, D = 10000, 3
x = torch.rand(N, D) # Point clouds in [0,1]<sup>3</sup>
y = torch.rand(N, D) # encoded as N-by-D arrays
# Re-indexing:
x_i = x[:,None,:] # shape (N,D) -> (N,1,D)
y_j = y[None,:,:] # shape (N,D) -> (1,N,D)
# Matrix of squared distances:
D_ij = ((x_i - y_j)**2).sum(dim=2) # (N,N,D) -> (N,N)
```

```
# Store arbitrary arrays on the GPU:
N, D = 10000, 3
x = torch.rand(N, D) # Point clouds in [0,1]<sup>3</sup>
y = torch.rand(N, D) # encoded as N-by-D arrays
# Re-indexing:
x_i = x[:,None,:] # shape (N,D) -> (N,1,D)
y_j = y[None,:,:] # shape (N,D) -> (1,N,D)
# Matrix of squared distances:
D_ij = ((x_i - y_j)**2).sum(dim=2) # (N,N,D) -> (N,N)
```

Memory overflow with 10k points? We should do better. [CL96]



11k triangles

871k triangles

What a GPU really looks like

Turing architecture.

PCI Espresa 3.0 Host histriace Oigethread Engine		
	L2 Cache	
erces erces erces erces erces erces	TTOM TOM TOM DIG DIG TOM	
	20 20 20 20 20 20 20 20 20 20 20 20 20 2	
GPC	are High-Speed Hub	GPC

What a GPU really looks like

The Admiralty and War Offices, built in 1884.



What a GPU really looks like

Inside view of a CUDA block.



• Blocks of ~200 threads.

- Blocks of ~200 threads.
- 4 layers of memory (Host, Device, Block, Thread).

- Blocks of ~200 threads.
- 4 layers of **memory** (Host, Device, Block, Thread).
- Shared Block buffer (~128 Kb) is 50x faster than Device mem.

- Blocks of ~200 threads.
- 4 layers of **memory** (Host, Device, Block, Thread).
- Shared Block buffer (~128 Kb) is 50x faster than Device mem.

- Blocks of ~200 threads.
- 4 layers of **memory** (Host, Device, Block, Thread).
- Shared Block buffer (~128 Kb) is 50x faster than Device mem.

PyTorch and TF variables are always stored in the **Device** mem.

- Blocks of ~200 threads.
- 4 layers of **memory** (Host, Device, Block, Thread).
- Shared Block buffer (~128 Kb) is 50x faster than Device mem.

PyTorch and TF variables are always stored in the **Device** mem. Explicit C++ implementations for:

- Convolutions.
- Fourier, wavelet transforms.

- Blocks of ~200 threads.
- 4 layers of **memory** (Host, Device, Block, Thread).
- Shared Block buffer (~128 Kb) is 50x faster than Device mem.

PyTorch and TF variables are always stored in the **Device** mem. Explicit C++ implementations for:

- Convolutions.
- Fourier, wavelet transforms.
- \implies We need to do the same!





Array Coefficients only



Array Coefficients only



Sparse matrix Coordinates + coeffs



Array Coefficients only **Symbolic LazyTensor** Formula + point clouds **Sparse matrix** Coordinates + coeffs



Coefficients only

Symbolic LazyTensor Formula + point clouds **Sparse matrix** Coordinates + coeffs

pip install pykeops

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(1000000, 3, requires_grad=True).cuda()
y = torch.rand(2000000, 3).cuda()
```

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(1000000, 3, requires_grad=True).cuda()
y = torch.rand(2000000, 3).cuda()
# Turn our Tensors into KeOps symbolic variables:
from pykeops.torch import LazyTensor
x i = LazyTensor(x[:,None,:]) # x i.shape = (1e6, 1, 3)
```

```
y_j = LazyTensor(y[None,:,:]) # y_j.shape = ( 1, 2e6,3)
```

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(1000000, 3, requires_grad=True).cuda()
y = torch.rand(2000000, 3).cuda()
# Turn our Tensors into KeOps symbolic variables:
from pykeops.torch import LazyTensor
x_i = LazyTensor(x[:,None,:]) # x_i.shape = (1e6, 1, 3)
```

```
y_j = LazyTensor(y[None,:,:]) # y_j.shape = ( 1, 2e6,3)
```

```
# Perform large-scale computations:
D_ij = ((x_i - y_j)**2).sum(dim=2) # (1e6,2e6,1)
```

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(1000000, 3, requires_grad=True).cuda()
y = torch.rand(2000000, 3).cuda()
# Turn our Tensors into KeOps symbolic variables:
from pykeops.torch import LazyTensor
x_i = LazyTensor(x[:,None,:]) # x_i.shape = (1e6, 1, 3)
y_j = LazyTensor(y[None,:,:]) # y_j.shape = (1, 2e6,3)
```

```
# Perform large-scale computations:
D_ij = ((x_i - y_j)**2).sum(dim=2) # (1e6,2e6,1)
K_ij = (- D_ij).exp() # (1e6,2e6,1)
```

```
# Large point clouds in [0,1]^3
import torch
x = torch.rand(1000000, 3, requires_grad=True).cuda()
y = torch.rand(2000000, 3).cuda()
# Turn our Tensors into KeOps symbolic variables:
from pykeops.torch import LazyTensor
x i = LazyTensor(x[:,None,:]) \# x i.shape = (1e6, 1, 3)
y_j = LazyTensor(y[None,:,:]) # y_j.shape = (1, 2e6,3)
# Perform large-scale computations:
D_{ij} = ((x_i - y_j) * 2).sum(dim=2) # (1e6,2e6,1)
K ij = (- D ij).exp()
                                 # (1e6,2e6,1)
# Reduction: symbolic LazyTensor -> genuine torch Tensor
a i = K ij.sum(dim=1) # (1e6, 1) torch.cuda.FloatTensor
```

```
# Large point clouds in [0,1]^3
import torch
x = torch.rand(1000000, 3, requires_grad=True).cuda()
y = torch.rand(2000000, 3).cuda()
# Turn our Tensors into KeOps symbolic variables:
from pykeops.torch import LazyTensor
x i = LazyTensor(x[:,None,:]) \# x i.shape = (1e6, 1, 3)
y_j = LazyTensor(y[None,:,:]) # y_j.shape = (1, 2e6,3)
# Perform large-scale computations:
D_{ij} = ((x_i - y_j) * 2).sum(dim=2) # (1e6,2e6,1)
K ij = (- D ij).exp()
                                    # (1e6,2e6,1)
# Reduction: symbolic LazyTensor -> genuine torch Tensor
a i = K ij.sum(dim=1) # (1e6, 1) torch.cuda.FloatTensor
g x = torch.autograd.grad((a_i ** 2).sum(), [x])
```

Scaling up to real data



Gaussian kernel product in 3D (RTX 2080 Ti GPU)

• 2 years of work with Joan Glaunès and Benjamin Charlier.

- 2 years of work with Joan Glaunès and Benjamin Charlier.
- Cross-platform: R, Matlab, NumPy and PyTorch.

- 2 years of work with Joan Glaunès and Benjamin Charlier.
- Cross-platform: R, Matlab, NumPy and PyTorch.
- Versatile: many operations, variables, reductions.

- 2 years of work with Joan Glaunès and Benjamin Charlier.
- Cross-platform: R, Matlab, NumPy and PyTorch.
- Versatile: many operations, variables, reductions.
- Efficient: O(N) memory, competitive runtimes.
The KeOps library

- 2 years of work with Joan Glaunès and Benjamin Charlier.
- Cross-platform: R, Matlab, NumPy and PyTorch.
- Versatile: many operations, variables, reductions.
- Efficient: O(N) memory, competitive runtimes.
- Powerful: autodiff, block-sparsity, etc.

The KeOps library

- 2 years of work with Joan Glaunès and Benjamin Charlier.
- Cross-platform: R, Matlab, NumPy and PyTorch.
- Versatile: many operations, variables, reductions.
- Efficient: O(N) memory, competitive runtimes.
- Powerful: autodiff, block-sparsity, etc.
- Transparent: interface with SciPy, GPytorch, etc.

The KeOps library

- 2 years of work with Joan Glaunès and Benjamin Charlier.
- Cross-platform: R, Matlab, NumPy and PyTorch.
- Versatile: many operations, variables, reductions.
- Efficient: O(N) memory, competitive runtimes.
- Powerful: autodiff, block-sparsity, etc.
- Transparent: interface with SciPy, GPytorch, etc.
- Scikit-learn-like documentation:

```
www.kernel-operations.io
```

Geometric loss functions

Life is easy when you have landmarks...



Anatomical landmarks from A morphometric approach for the analysis of body shape in bluefin tuna, Addis et al., 2009.

Life is easy when you have landmarks...



Anatomical landmarks from A morphometric approach for the analysis of body shape in bluefin tuna, Addis et al., 2009.

Unfortunately, medical data is often unlabeled [EPW⁺11]



Surface meshes



Segmentation masks

Let's enforce sampling invariance:

$$A \longrightarrow \alpha = \sum_{i=1}^{N} \alpha_i \delta_{\mathbf{x}_i}, \qquad B \longrightarrow \beta = \sum_{j=1}^{M} \beta_j \delta_{\mathbf{y}_j}.$$

Let's enforce sampling invariance:

$$\mathsf{A} \ \longrightarrow \ \alpha \ = \ \sum_{i=1}^{\mathsf{N}} \alpha_i \delta_{\mathsf{x}_i} \,, \qquad \mathsf{B} \ \longrightarrow \ \beta \ = \ \sum_{j=1}^{\mathsf{M}} \beta_j \delta_{y_j} \,.$$



Let's enforce sampling invariance:

$$\mathsf{A} \ \longrightarrow \ \alpha \ = \ \sum_{i=1}^{\mathsf{N}} \alpha_i \delta_{\mathsf{x}_i} \,, \qquad \mathsf{B} \ \longrightarrow \ \beta \ = \ \sum_{j=1}^{\mathsf{M}} \beta_j \delta_{\mathsf{y}_j} \,.$$







$$\alpha = \sum_{i=1}^{N} \alpha_i \delta_{\mathbf{x}_i}, \quad \beta = \sum_{j=1}^{M} \beta_j \delta_{\mathbf{y}_j}.$$



$$\alpha = \sum_{i=1}^{N} \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^{M} \beta_j \delta_{y_j}.$$
$$\sum_{i=1}^{N} \alpha_i = 1 = \sum_{j=1}^{M} \beta_j$$

(



$$\alpha = \sum_{i=1}^{N} \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^{M} \beta_j \delta_{y_j}.$$
$$\sum_{i=1}^{N} \alpha_i = 1 = \sum_{j=1}^{M} \beta_j$$
Display $v_i = -\frac{1}{\alpha_i} \nabla_{x_i} \text{Loss}(\alpha, \beta).$



$$\alpha = \sum_{i=1}^{N} \alpha_i \delta_{\mathbf{x}_i}, \quad \beta = \sum_{j=1}^{M} \beta_j \delta_{\mathbf{y}_j}.$$
$$\sum_{i=1}^{N} \alpha_i = 1 = \sum_{j=1}^{M} \beta_j$$
Display $\nu_i = -\frac{1}{\alpha_i} \nabla_{\mathbf{x}_i} \text{Loss}(\alpha, \beta).$

Seamless extensions to:

- $\sum_{i} \alpha_{i} \neq \sum_{j} \beta_{j}$, outliers [CPSV18],
- curves and surfaces [KCC17],
- variable weights α_i .

Simple loss functions

• Chamfer distance \simeq soft-Hausdorff: Projection-based \longrightarrow Degenerate gradients.

Simple loss functions

- Chamfer distance \simeq soft-Hausdorff: Projection-based \longrightarrow Degenerate gradients.
- Kernel distance \simeq Blurred SSD, convolution-based: Loss $(\alpha, \beta) = \frac{1}{2} ||g \star (\alpha - \beta)||_{L^2(\mathbb{R}^D)}^2 = \frac{1}{2} \langle \alpha - \beta, k \star (\alpha - \beta) \rangle$ where $k = (g \circ (x \mapsto -x)) \star g$.

Simple loss functions

- Chamfer distance \simeq soft-Hausdorff: Projection-based \longrightarrow Degenerate gradients.
- Kernel distance \simeq Blurred SSD, convolution-based: Loss $(\alpha, \beta) = \frac{1}{2} ||g \star (\alpha - \beta)||^2_{L^2(\mathbb{R}^D)} = \frac{1}{2} \langle \alpha - \beta, k \star (\alpha - \beta) \rangle$ where $k = (g \circ (x \mapsto -x)) \star g$.
- Example: the Energy Distance, k(x, y) = -||x y||:

$$Loss(\alpha, \beta) = \sum_{i} \sum_{j} \alpha_{i} \beta_{j} \|\mathbf{x}_{i} - \mathbf{y}_{j}\| \\ - \frac{1}{2} \sum_{i} \sum_{j} \alpha_{i} \alpha_{j} \|\mathbf{x}_{i} - \mathbf{x}_{j}\| - \frac{1}{2} \sum_{i} \sum_{j} \beta_{i} \beta_{j} \|\mathbf{y}_{i} - \mathbf{y}_{j}\|.$$



$$t = .00$$



$$t = .25$$





$$t=1.00$$



$$t = 5.00$$



We need **clean gradients**, without artifacts.

We need **clean gradients**, without artifacts.

We need clean gradients, without artifacts.



We need clean gradients, without artifacts.



We need **clean gradients**, without artifacts.



We need clean gradients, without artifacts.



$$OT(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2N} \sum_{i=1}^{N} |\mathbf{x}_i - \mathbf{y}_{\sigma^*(i)}|^2$$

We need clean gradients, without artifacts.



$$OT(\alpha, \beta) = \frac{1}{2N} \sum_{i=1}^{N} |\mathbf{x}_{i} - y_{\sigma^{*}(i)}|^{2} = \min_{\sigma \in S_{N}} \frac{1}{2N} \sum_{i=1}^{N} |\mathbf{x}_{i} - y_{\sigma(i)}|^{2}$$

Optimal transport generalizes sorting to ${\sf D}>1$



Minimize over N-by-M matrices (transport plans) π :

$$OT(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \min_{\pi} \underbrace{\sum_{i,j} \pi_{i,j} \cdot \frac{1}{2} |\mathbf{x}_i - \mathbf{y}_j|^2}_{\text{transport cost}}$$



subject to $\pi_{i,j} \ge 0$, $\sum_{j} \pi_{i,j} = \alpha_{i}, \quad \sum_{i} \pi_{i,j} = \beta_{j}.$



$$t = .00$$



$$t = .25$$



$$t = .50$$



$$t = 1.00$$
Wasserstein gradients are homogeneous



$$t = 5.00$$

Wasserstein gradients are homogeneous



The Wasserstein loss $\mathsf{OT}(\alpha,\beta)$ is:

• Symmetric: $OT(\alpha, \beta) = OT(\beta, \alpha)$.

The Wasserstein loss $OT(\alpha, \beta)$ is:

- Symmetric: $OT(\alpha, \beta) = OT(\beta)$
- Positive:

 $\operatorname{OT}(\boldsymbol{\alpha},\boldsymbol{\beta}) = \operatorname{OT}(\boldsymbol{\beta},\boldsymbol{\alpha}).$ $\operatorname{OT}(\boldsymbol{\alpha},\boldsymbol{\beta}) \geqslant 0.$

- Symmetric: $OT(\alpha, \beta) = OT(\beta, \alpha)$.
- Positive: $OT(\alpha, \beta) \ge 0$.
- Definite:

$$OT(\alpha, \beta) \ge 0.$$
$$OT(\alpha, \beta) = 0 \iff \alpha = \beta.$$

- Symmetric: $OT(\alpha, \beta) = OT(\beta, \alpha)$.
- Positive: $OT(\boldsymbol{\alpha}, \boldsymbol{\beta}) \ge 0$.
- Definite: $\operatorname{OT}(\alpha,\beta) = 0 \Longleftrightarrow \alpha = \beta$.
- Translation-aware: $OT(\alpha, Translate_{\vec{v}}(\alpha)) = \frac{1}{2} \|\vec{v}\|^2$.

- Symmetric: $OT(\alpha, \beta) = OT(\beta, \alpha)$.
- Positive: $OT(\boldsymbol{\alpha}, \boldsymbol{\beta}) \ge 0$.
- Definite: $OT(\alpha, \beta) = 0 \iff \alpha = \beta$.
- Translation-aware: $OT(\alpha, Translate_{\vec{v}}(\alpha)) = \frac{1}{2} \|\vec{v}\|^2$.
- More generally, OT retrieves the unique gradient of a convex function T = ∇φ that maps α onto β:
 - $\begin{array}{ll} \text{In dimension 1,} & (\textbf{x}_i \textbf{x}_j) \cdot (\textbf{y}_{\sigma(i)} \textbf{y}_{\sigma(j)}) & \geqslant 0 \\ \text{In dimension D,} & \langle \textbf{x}_i \textbf{x}_j \ , \ \textbf{T}(\textbf{x}_i) \textbf{T}(\textbf{x}_j) \rangle_{\mathbb{R}^D} & \geqslant 0 \ . \end{array}$

- Symmetric: $OT(\alpha, \beta) = OT(\beta, \alpha)$.
- Positive: $OT(\boldsymbol{\alpha}, \boldsymbol{\beta}) \ge 0$.
- Definite: $OT(\alpha, \beta) = 0 \iff \alpha = \beta$.
- Translation-aware: $OT(\alpha, Translate_{\vec{v}}(\alpha)) = \frac{1}{2} \|\vec{v}\|^2$.
- More generally, OT retrieves the unique gradient of a convex function T = ∇φ that maps α onto β:
 - $\begin{array}{ll} \text{In dimension 1,} & (\textbf{x}_i \textbf{x}_j) \cdot (\textbf{y}_{\sigma(i)} \textbf{y}_{\sigma(j)}) & \geqslant 0 \\ \text{In dimension D,} & \langle \textbf{x}_i \textbf{x}_j \ , \ \textbf{T}(\textbf{x}_i) \textbf{T}(\textbf{x}_j) \rangle_{\mathbb{R}^D} & \geqslant 0 \ . \end{array}$

The Wasserstein loss $OT(\alpha, \beta)$ is:

- Symmetric: $OT(\alpha, \beta) = OT(\beta, \alpha)$.
- Positive: $OT(\boldsymbol{\alpha}, \boldsymbol{\beta}) \ge 0$.
- Definite: $OT(\alpha, \beta) = 0 \iff \alpha = \beta$.
- Translation-aware: $OT(\alpha, Translate_{\vec{v}}(\alpha)) = \frac{1}{2} \|\vec{v}\|^2$.
- More generally, OT retrieves the unique gradient of a convex function T = ∇φ that maps α onto β:
 - $\begin{array}{ll} \text{In dimension 1,} & (\textbf{x}_i \textbf{x}_j) \cdot (\textbf{y}_{\sigma(i)} \textbf{y}_{\sigma(j)}) & \geqslant 0 \\ \text{In dimension D,} & \langle \textbf{x}_i \textbf{x}_j \ , \ \textbf{T}(\textbf{x}_i) \textbf{T}(\textbf{x}_j) \rangle_{\mathbb{R}^D} & \geqslant 0 \ . \end{array}$

 \implies Appealing generalization of an **increasing mapping**.

Key dates:

• [Kan42]: Dual problem.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.
- [GRL+98, CR00]: Robust Point Matching = Sinkhorn as a loss.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.
- [GRL+98, CR00]: Robust Point Matching = Sinkhorn as a loss.
- [Cut13]: Start of the GPU era.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.
- [GRL+98, CR00]: Robust Point Matching = Sinkhorn as a loss.
- [Cut13]: Start of the GPU era.
- [Mér11, Lév15, Sch19]: Multiscale CPU solvers in O(N log N).

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.
- [GRL+98, CR00]: Robust Point Matching = Sinkhorn as a loss.
- [Cut13]: Start of the GPU era.
- [Mér11, Lév15, Sch19]: Multiscale CPU solvers in O(N log N).
- Today: Multiscale Sinkhorn algorithm, on the GPU.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.
- [GRL+98, CR00]: Robust Point Matching = Sinkhorn as a loss.
- [Cut13]: Start of the GPU era.
- [Mér11, Lév15, Sch19]: Multiscale CPU solvers in O(N log N).
- Today: Multiscale Sinkhorn algorithm, on the GPU.

- [Kan42]: Dual problem.
- [Kuh55]: Hungarian method in $O(N^3)$.
- [Ber79]: Auction algorithm in $O(N^2)$.
- [KY94]: **SoftAssign** = Sinkhorn + annealing, in $O(N^2)$.
- [GRL+98, CR00]: Robust Point Matching = Sinkhorn as a loss.
- [Cut13]: Start of the GPU era.
- [Mér11, Lév15, Sch19]: Multiscale CPU solvers in O(N log N).
- Today: Multiscale Sinkhorn algorithm, on the GPU.

 \implies Generalized **QuickSort** algorithm.



















+ $O(N \log N)$ instead of $O(N^2)$.

- + $O(N \log N)$ instead of $O(N^2)$.
- +~ Converge in $\,\leqslant\,$ 10 iterations.

- + $O(N \log N)$ instead of $O(N^2)$.
- +~ Converge in $\,\leqslant\,$ 10 iterations.
- + Fully symmetric.

- + $O(N \log N)$ instead of $O(N^2)$.
- +~ Converge in $\,\leqslant\,$ 10 iterations.
- + Fully symmetric.
- +~ Positive and definite $[\text{FSV}^+18].$

- + $O(N \log N)$ instead of $O(N^2)$.
- +~ Converge in $\,\leqslant\,$ 10 iterations.
- + Fully symmetric.
- +~ Positive and definite $[\text{FSV}^+18].$
- Much harder to implement.

- + $O(N \log N)$ instead of $O(N^2)$.
- +~ Converge in $\,\leqslant\,$ 10 iterations.
- + Fully symmetric.
- +~ Positive and definite $[\text{FSV}^+18].$
- Much harder to implement.

- + $O(N \log N)$ instead of $O(N^2)$.
- +~ Converge in $\,\leqslant\,$ 10 iterations.
- + Fully symmetric.
- + Positive and definite [FSV⁺18].
- Much harder to implement.

Geometric Loss functions for PyTorch

Our website: www.kernel-operations.io/geomloss

\implies pip install geomloss \Leftarrow
\implies pip install geomloss \Leftarrow

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(100000, 3, requires_grad=True).cuda()
y = torch.rand(200000, 3).cuda()
```

 \implies pip install geomloss \Leftarrow

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(100000, 3, requires_grad=True).cuda()
y = torch.rand(200000, 3).cuda()
# Define a Wasserstein loss between sampled measures
from geomloss import SamplesLoss # See also ImagesLoss...
```

```
loss = SamplesLoss(loss="sinkhorn", p=2, blur=.05)
```

 \implies pip install geomloss \Leftarrow

```
# Large point clouds in [0,1]<sup>3</sup>
import torch
x = torch.rand(100000, 3, requires_grad=True).cuda()
y = torch.rand(200000, 3).cuda()
# Define a Wasserstein loss between sampled measures
from geomloss import SamplesLoss # See also ImagesLoss...
loss = SamplesLoss(loss="sinkhorn", p=2, blur=.05)
```

L = loss(x, y) # By default, use constant weights

 \Rightarrow pip install geomloss \Leftarrow

```
# Large point clouds in [0,1]^3
import torch
x = torch.rand(100000, 3, requires_grad=True).cuda()
v = torch.rand(200000, 3).cuda()
# Define a Wasserstein loss between sampled measures
from geomloss import SamplesLoss # See also ImagesLoss...
loss = SamplesLoss(loss="sinkhorn", p=2, blur=.05)
L = loss(x, y) # By default, use constant weights
# GeomLoss supports autograd, batch processing, etc.
g_x, = torch.autograd.grad(L, [x])
```

Scaling up optimal transport

Precision controlled by the ratio diameter

With a precision of 1%, on a modern gaming GPU:



blur

10k points in 30-50ms

100k points in 100-200ms

Affordable geometric barycenters

Barycenter
$$\alpha^* = \arg \min_{\alpha} \sum_{i=1}^N \lambda_i \operatorname{Loss}(\alpha, \beta_i)$$
.

Affordable geometric barycenters



Affordable geometric barycenters



A robust loss function



A robust loss function



A robust loss function

















Iteration 2



















Conclusion

- Symbolic LazyTensors are key to performances:
 - \longrightarrow KeOps, x30 speed-up vs PyTorch and TF.

- Symbolic LazyTensors are key to performances:
 - \longrightarrow KeOps, x30 speed-up vs PyTorch and TF.
- Optimal Transport = generalized sorting:
 - \longrightarrow Geometric gradients.
 - \longrightarrow Super-fast $O(N \log N)$ solvers.

- Symbolic LazyTensors are key to performances:
 - \longrightarrow KeOps, x30 speed-up vs PyTorch and TF.
- Optimal Transport = generalized sorting:
 - \longrightarrow Geometric gradients.
 - \longrightarrow Super-fast $O(N \log N)$ solvers.

• Symbolic LazyTensors are key to performances:

 \longrightarrow KeOps, x30 speed-up vs PyTorch and TF.

- Optimal Transport = generalized sorting:
 - \longrightarrow Geometric gradients.
 - \longrightarrow Super-fast $O(N \log N)$ solvers.

 \implies www.kernel-operations.io \Leftarrow

- Symbolic LazyTensors are key to performances:
 - \longrightarrow KeOps, x30 speed-up vs PyTorch and TF.
- Optimal Transport = generalized sorting:
 - \longrightarrow Geometric gradients.
 - \longrightarrow Super-fast $O(N \log N)$ solvers.

\implies www.kernel-operations.io \Leftarrow

Main reference: my PhD thesis, available in January.

Topology-aware registration is a hard, non-convex problem. OT is good enough as a **loss**, but is *not* the panacea. **Topology-aware** registration is a hard, non-convex problem. OT is good enough as a **loss**, but is *not* the panacea.

 \implies We need robust, data-driven shape metrics.

Topology-aware registration is a hard, non-convex problem. OT is good enough as a **loss**, but is *not* the panacea.

\implies We need robust, data-driven shape metrics. (3^{\rm rd} half)
My work so far: theory, implementation and proofs of concept





My work so far: theory, implementation and proofs of concept



Theorems, libraries, documentation... But no real validation :- (

Thank you for your attention.

Let's have a chat!

First setting: processing of point clouds



- + φ is \mathbf{rigid} or affine
- Occlusions
- Outliers

From the documentation of the Point Cloud Library.

Second setting: medical imaging



From Marc Niethammer's Quicksilver slides.

- φ is a spline or a **diffeomorphism**
- Ill-posed problem
- Some occlusions



Wasserstein Auto-Encoders, Tolstikhin et al., 2018.

- + φ is a neural network
- Very weak regularization
- High-dimensional space



- + φ is a neural network
- Very weak regularization
- High-dimensional space

Wasserstein Auto-Encoders, Tolstikhin et al., 2018.

Which **Loss** function should we use?

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

look for $\theta^* = \arg \min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

$$\mathsf{Loss}(\alpha,\beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

look for $\theta^* = \arg\min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

•
$$B = \{ \|f\|_{\infty} \leq 1 \} \implies \text{Loss} = \text{TV norm:}$$

- zero geometry
- too many test functions

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

look for $\theta^* = \arg\min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

•
$$B = \{ \|f\|_{\infty} \leq 1 \} \implies \text{Loss} = \text{TV norm:}$$

- zero geometry
- too many test functions
- $B = \{ \|f\|_2^2 + \|\nabla f\|_2^2 + \dots \leq 1 \} \Longrightarrow$ Loss = kernel norm:
 - may saturate at infinity
 - screening artifacts

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

look for $\theta^* = \arg \min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

• $B = \{ f \text{ is } 1 \text{-Lipschitz} \} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0):$

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

look for $\theta^* = \arg \min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

- $B = \{ f \text{ is } 1 \text{-Lipschitz} \} \implies \text{Loss} = \text{Wasserstein-1 (OT_0):}$
 - + S $_{\varepsilon}$ is nearly as efficient as a closed formula

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

look for $\theta^* = \arg\min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

- $B = \{ f \text{ is } 1 \text{-Lipschitz} \} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0):$
 - + S $_{\varepsilon}$ is nearly as efficient as a closed formula
 - relevant in low dimensions
 - useless in ($\mathbb{R}^{512 \times 512}, \|\cdot\|_2$): the ground cost makes no sense

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

ook for $\theta^* = \arg \min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

- $B = \{f \text{ is } 1 \text{-Lipschitz} \} \Longrightarrow \text{ Loss} = \text{Wasserstein-1 (OT_0):}$
 - + S $_{\varepsilon}$ is nearly as efficient as a closed formula
 - relevant in low dimensions
 - useless in $(\mathbb{R}^{512 \times 512}, \|\cdot\|_2)$: the ground cost makes no sense
- $B \simeq \{ f \text{ is 1-Lipschitz } \} \bigcap \{ f \text{ is a CNN } \}$ $\implies \text{Loss} = \text{Wasserstein GAN }:$

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

ook for $\theta^* = \arg \min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

- $B = \{f \text{ is } 1 \text{-Lipschitz} \} \Longrightarrow \text{ Loss} = \text{Wasserstein-1 (OT_0):}$
 - + S $_{\varepsilon}$ is nearly as efficient as a closed formula
 - relevant in low dimensions
 - useless in $(\mathbb{R}^{512 \times 512}, \|\cdot\|_2)$: the ground cost makes no sense
- $B \simeq \{ f \text{ is } 1 \text{-Lipschitz} \} \bigcap \{ f \text{ is a CNN} \}$

 \implies Loss = Wasserstein GAN :

• use perceptually sensible test functions

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

ook for $\theta^* = \arg \min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

- $B = \{f \text{ is } 1 \text{-Lipschitz} \} \Longrightarrow \text{ Loss} = \text{Wasserstein-1 (OT_0):}$
 - + S $_{\varepsilon}$ is nearly as efficient as a closed formula
 - relevant in low dimensions
 - useless in ($\mathbb{R}^{512\times512},\|\cdot\|_2$): the ground cost makes no sense
- $B \simeq \{f \text{ is } 1\text{-Lipschitz }\} \bigcap \{f \text{ is a CNN }\}$

 \implies Loss = Wasserstein GAN :

- use perceptually sensible test functions
- no simple formula: use gradient ascent

$$Loss(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

ook for $\theta^* = \arg\min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$

- $B = \{f \text{ is } 1 \text{-Lipschitz} \} \Longrightarrow \text{ Loss} = \text{Wasserstein-1 (OT_0):}$
 - + S $_{\varepsilon}$ is nearly as efficient as a closed formula
 - relevant in low dimensions
 - useless in ($\mathbb{R}^{512\times512},\|\cdot\|_2$): the ground cost makes no sense
- $B \simeq \{ f \text{ is } 1 \text{-Lipschitz} \} \bigcap \{ f \text{ is a CNN} \}$

 \implies Loss = Wasserstein GAN :

- use **perceptually sensible** test functions
- no simple formula: use gradient ascent
- can we provide relevant **insights** to the ML community?

Our papers:

Global divergences between measures: from Hausdorff distance to
 Optimal Transport, F., Trouvé, 2018

Our papers:

- Global divergences between measures: from Hausdorff distance to
 Optimal Transport, F., Trouvé, 2018
- Sinkhorn entropies and divergences,
 - F., Séjourné, Vialard, Amari, Trouvé, Peyré, 2018

Our papers:

- Global divergences between measures: from Hausdorff distance to
 Optimal Transport, F., Trouvé, 2018
- Sinkhorn entropies and divergences,
 F., Séjourné, Vialard, Amari, Trouvé, Peyré, 2018
- Optimal Transport for diffeomorphic registration, F., Charlier, Vialard, Peyré, 2017

References i



John Ashburner.

A fast diffeomorphic image registration algorithm.

Neuroimage, 38(1):95–113, 2007.

Dimitri P Bertsekas.

A distributed algorithm for the assignment problem.

Lab. for Information and Decision Systems Working Paper, M.I.T., Cambridge, MA, 1979.



Y. Brenier.

Polar factorization and monotone rearrangement of vector-valued functions.

Comm. Pure Appl. Math., 44(4):375-417, 1991.

Brian Curless and Marc Levoy.

A volumetric method for building complex models from range images.

In Proceedings of the 23rd annual conference on Computer graphics and interactive techniques, pages 303–312. ACM, 1996.

Christophe Chnafa, Simon Mendez, and Franck Nicoud.
 Image-based large-eddy simulation in a realistic left heart.
 Computers & Fluids, 94:173–187, 2014.

 Lénaïc Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard.
 Unbalanced optimal transport: Dynamic and kantorovich formulations.

Journal of Functional Analysis, 274(11):3090–3123, 2018.

Haili Chui and Anand Rangarajan.

A new algorithm for non-rigid point matching.

In Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on, volume 2, pages 44–51. IEEE, 2000.

References iv

Adam Conner-Simons and Rachel Gordon. **Using ai to predict breast cancer and personalize care.** http://news.mit.edu/2019/ using-ai-predict-breast-cancer-and-personalize-c 2019. MIT CSAIL.

- Marco Cuturi.

Sinkhorn distances: Lightspeed computation of optimal transport.

In Advances in Neural Information Processing Systems, pages 2292–2300, 2013.

References v

Olivier Ecabert, Jochen Peters, Matthew J Walker, Thomas Ivanc, Cristian Lorenz, Jens von Berg, Jonathan Lessick, Mani Vembar, and Jürgen Weese.

Segmentation of the heart and great vessels in CT images using a model-based adaptation framework. Medical image analysis, 15(6):863–876, 2011.

 Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-Ichi Amari, Alain Trouvé, and Gabriel Peyré.
 Interpolating between optimal transport and MMD using Sinkhorn divergences.

arXiv preprint arXiv:1810.08278, 2018.



Joan Glaunes.

Transport par difféomorphismes de points, de mesures et de courants pour la comparaison de formes et l'anatomie numérique.

These de sciences, Université Paris, 13, 2005.

Steven Gold, Anand Rangarajan, Chien-Ping Lu, Suguna Pappu, and Eric Mjolsness.

New algorithms for 2d and 3d point matching: Pose estimation and correspondence.

Pattern recognition, 31(8):1019–1031, 1998.

Leonid V Kantorovich.

On the translocation of masses.

In Dokl. Akad. Nauk. USSR (NS), volume 37, pages 199–201, 1942.

 Irene Kaltenmark, Benjamin Charlier, and Nicolas Charon.
 A general framework for curve and surface comparison and registration with oriented varifolds.

In Computer Vision and Pattern Recognition (CVPR), 2017.

] Harold W Kuhn.

The Hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.

🔋 Jeffrey J Kosowsky and Alan L Yuille.

The invisible hand algorithm: Solving the assignment problem with statistical physics.

Neural networks, 7(3):477-490, 1994.



Bruno Lévy.

A numerical algorithm for l2 semi-discrete optimal transport in 3d.

ESAIM: Mathematical Modelling and Numerical Analysis, 49(6):1693–1715, 2015.

References ix

Christian Ledig, Andreas Schuh, Ricardo Guerrero, Rolf A Heckemann, and Daniel Rueckert. Structural brain imaging in Alzheimer's disease and mild cognitive impairment: biomarker analysis and shared morphometry database.

Scientific reports, 8(1):11258, 2018.

Stéphane Mallat.

Understanding deep convolutional networks.

Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 374(2065):20150203, 2016.



Quentin Mérigot.

A multiscale approach to optimal transport.

In *Computer Graphics Forum*, volume 30, pages 1583–1592. Wiley Online Library, 2011.

Yaroslav Nikulin and Roman Novak.

Exploring the neural algorithm of artistic style.

arXiv preprint arXiv:1602.07188, 2016.

Moses Olafenwa.

Object detection with 10 lines of code. https://towardsdatascience.com/ object-detection-with-10-lines-of-code-d6cb4d861 2018.

Towards Data Science.

Maurice Peemen, Bart Mesman, and Henk Corporaal.
Speed sign detection and recognition by convolutional neural networks.

In Proceedings of the 8th international automotive congress, pages 162–170. sn, 2011.



Ptrump16.

Irm picture.

https://commons.wikimedia.org/w/index.php? curid=64157788,2019. CC BY-SA 4.0. Olaf Ronneberger, Philipp Fischer, and Thomas Brox.
 U-net: Convolutional networks for biomedical image segmentation.

In International Conference on Medical image computing and computer-assisted intervention, pages 234–241. Springer, 2015.

Bernhard Schmitzer.
 Stabilized sparse scaling algorithms for entropy regularized transport problems.

SIAM Journal on Scientific Computing, 41(3):A1443–A1481, 2019.

 Donglai Wei, Bolei Zhou, Antonio Torralba, and William T Freeman.
 mNeuron: A Matlab plugin to visualize neurons from deep models.
 Massachusetts Institute of Technology, 2017.