

# Fast geometric learning with symbolic matrices

---

Jean Feydy  
Imperial College London

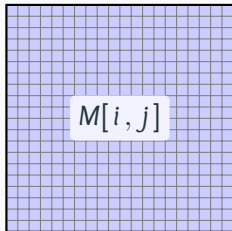
CogSys seminar, DTU Compute, online — January 2021.

Joint work with B. Charlier, J. Glaunès (numerical foundations),  
T. Séjourné, F.-X. Vialard, G. Peyré (optimal transport theory),  
P. Roussillon, P. Gori, A. Trouvé (applications to computational anatomy),  
F. Sverrisson, B. E. Correia, M. Bronstein (applications to protein sciences).

Symbolic matrices?

---

# Machine learning libraries represent most objects as tensors



## Dense matrix

Coefficients only

**Dense** matrices – large, contiguous **arrays** of numbers:

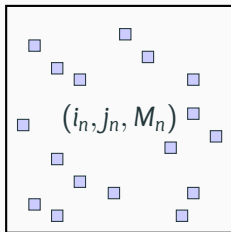
- + **Convenient** and well supported.
- Heavy load on the **memories** of our GPUs, with **time-consuming transfers** that take place between compute units.

# Machine learning libraries represent most objects as tensors



Dense matrix

Coefficients only



Sparse matrix

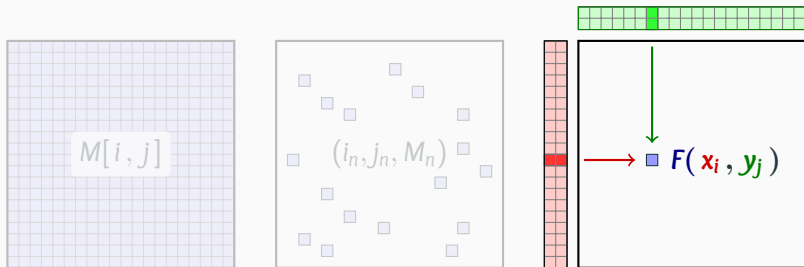
Coordinates + coeffs

**Sparse matrices** – tensors that have **few non-zero entries**:

- + Represent **large tensors** with a small memory footprint.
- Outside of **graph** processing, few objects are **sparse enough** to really benefit from this representation.



# Machine learning libraries represent most objects as tensors



Dense matrix

Coefficients only

Sparse matrix

Coordinates + coeffs

Symbolic matrix

Formula + data

**Distance** and **kernel** matrices, **point** convolutions, **attention** layers:

- + **Linear** memory usage: no more **memory** overflows.
- + We can optimize the use of registers for a  $\times 10 - \times 100$  **speed-up** vs. a standard PyTorch GPU baseline.

# We provide support for this “new abstraction” on the GPU

**Our library** comes with all the perks of a deep learning toolbox:

- + Transparent **array-like** interface.
- + Full support for automatic **differentiation**.
- + Comprehensive collection of **tutorials**, available online.

Under the hood: combines an optimized **C++** engine with high-level binders for **PyTorch**, **NumPy**, Matlab and R (thanks to Ghislain Durif).  
(We welcome **contributors** for JAX, Julia and other frameworks!)

To get started:

```
⇒ pip install pykeops ←  
www.kernel-operations.io
```

## A first example: efficient nearest neighbor search in dimension 50

Create large point clouds using **standard PyTorch syntax**:

```
import torch
N, M, D = 10**6, 10**6, 50
x = torch.rand(N, 1, D).cuda() # (1M, 1, 50) array
y = torch.rand(1, M, D).cuda() # (1, 1M, 50) array
```

Turn **dense** arrays into **symbolic** matrices:

```
from pykeops.torch import LazyTensor
x_i, y_j = LazyTensor(x), LazyTensor(y)
```

Create a large **symbolic matrix** of squared distances:

```
D_ij = ((x_i - y_j)**2).sum(dim=2) # (1M, 1M) symbolic
```

Use an `.argmin()` **reduction** to perform a nearest neighbor query:

```
indices_i = D_ij.argmax(dim=1) # -> standard torch tensor
```

## The KeOps library combines performance with flexibility

Script of the previous slide = efficient nearest neighbor query,  
**on par** with the bruteforce CUDA scheme of the **FAISS** library...

And can be used with **any metric!**

```
D_ij = ((x_i - x_j) ** 2).sum(dim=2)      # Euclidean
M_ij = (x_i - x_j).abs().sum(dim=2)     # Manhattan
C_ij = 1 - (x_i | x_j)                  # Cosine
H_ij = D_ij / (x_i[...,0] * x_j[...,0]) # Hyperbolic
```

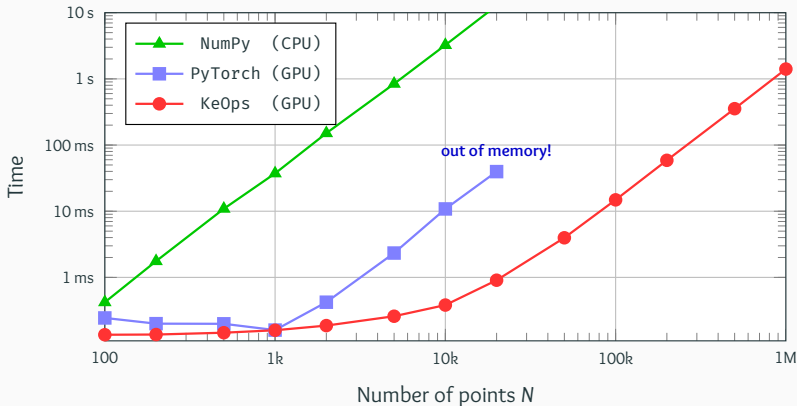
KeOps supports arbitrary **formulas** and **variables** with:

- **Reductions:** sum, log-sum-exp, K-min, matrix-vector product, etc.
- **Operations:** +, ×, sqrt, exp, neural networks, etc.
- **Advanced schemes:** batch processing, block sparsity, etc.
- **Automatic differentiation:** seamless integration with PyTorch.

# KeOps lets users work with millions of points at a time

Benchmark of a matrix-vector product with a N-by-N Gaussian kernel matrix between 3D point clouds.

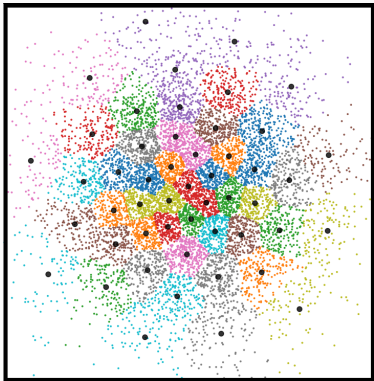
We run NumPy, PyTorch and KeOps on a RTX 2080 Ti GPU.



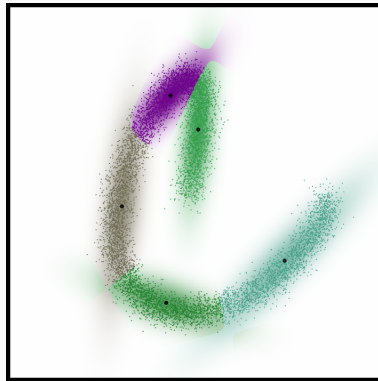
# Applications

---

# KeOps is a good fit for machine learning research



K-Means.

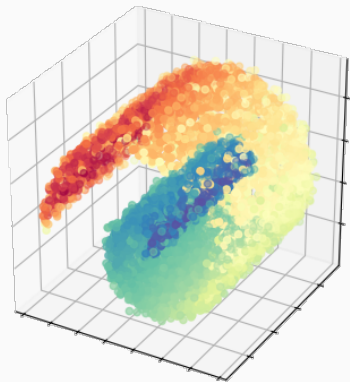


Gaussian Mixture Model.

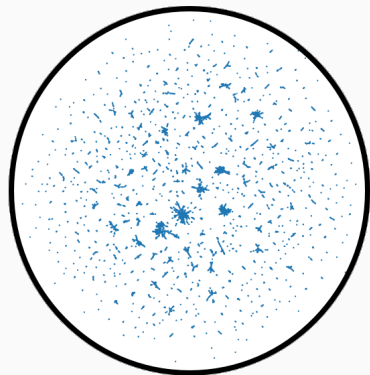
**Use any kernel, metric or formula you like!**

⇒ More tutorials coming up soon.

# KeOps is a good fit for machine learning research



Spectral analysis.



UMAP in hyperbolic space.

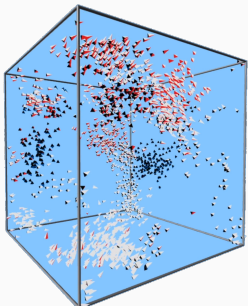
**Use any kernel, metric or formula you like!**

⇒ More tutorials coming up soon.

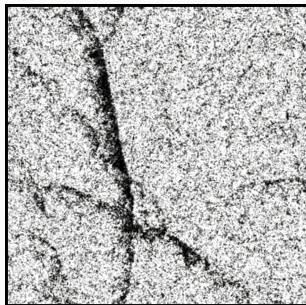


# KeOps lets you focus on your models, results and theorems

Some applications to **dynamical systems** [DM08, DFMAT17] and **statistics** [CDF19] with A. Diez, G. Clarté and P. Degond:



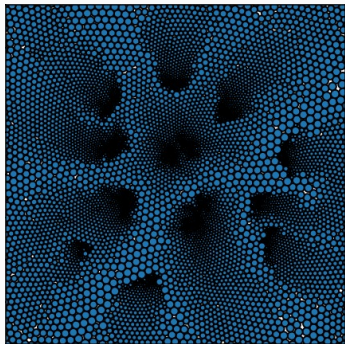
3D Vicsek model with orientation,  
interactive demo with 2k **flyers**.



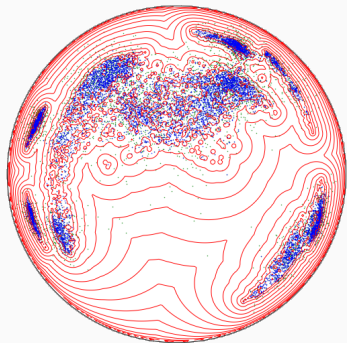
2D Vicsek model on the torus,  
in real-time with 100k **swimmers**.

# KeOps lets you focus on your models, results and theorems

⇒ Scale up to millions/billions of agents with Python scripts.

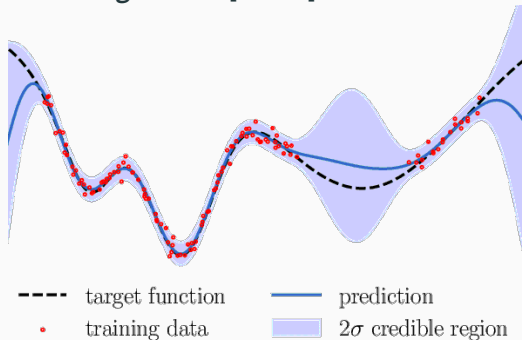


**Packing** problem in 2D  
with 10k repulsive balls.



Collective Monte Carlo **sampling**  
on the hyperbolic Poincaré disk.

A standard tool for regression [Lec18]:



Under the hood, solve a **kernel linear system**:

$$(\lambda \text{Id} + K_{xx}) a = b \quad \text{i.e.} \quad a \leftarrow (\lambda \text{Id} + K_{xx})^{-1} b$$

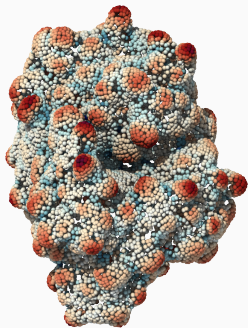
where  $\lambda \geq 0$  and  $(K_{xx})_{i,j} = k(x_i, x_j)$  is a positive definite matrix.

## KeOps symbolic tensors:

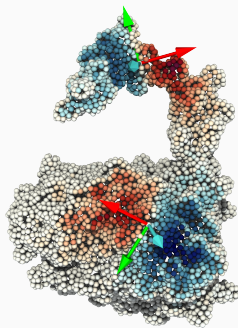
- Can be fed to **standard solvers**: SciPy, GPytorch, etc.
- GPytorch on the 3DRoad dataset ( $N = 278k, D = 3$ ):  
7h with 8 GPUs → 15mn with 1 GPU.
- Provide a **fast backend for research codes**: see e.g.  
*Kernel methods through the roof: handling **billions of points** efficiently*, by G. Meanti, L. Carratino, L. Rosasco, A. Rudi (2020).

Data-driven methods on **point clouds** and **proteins** [SFCB20]:

- + **Fast K-NN search**: local interactions.
- + **Fast N-by-N computations**: global interactions.
- + Heterogeneous **batches**, Octree-like acceleration.



Curvatures at all scales.



Quasi-geodesic convolutions.

**Fast, scalable and robust  
optimal transport solvers**

---

## How should we solve the OT problem?

Key dates for discrete optimal transport with  $N$  points:

- [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<sup>+</sup>98, CR00]: **Robust Point Matching** = Sinkhorn as a loss.
- [Cut13]: Start of the **GPU era**.
- [Mér11, Lév15, Sch19]: **Multiscale** solvers in  $O(N \log N)$ .
- Today: **Multiscale Sinkhorn algorithm, on the GPU**.

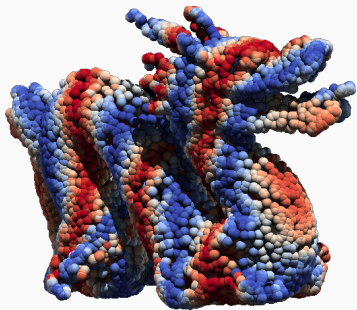
⇒ Generalized **QuickSort** algorithm.

# Scaling up optimal transport to anatomical data

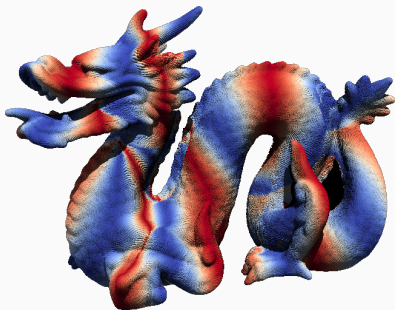
Progresses of the last decade add up to a  $\times 100$  -  $\times 1000$  acceleration:

Sinkhorn GPU  $\xrightarrow{\times 10}$  + KeOps  $\xrightarrow{\times 10}$  + Annealing  $\xrightarrow{\times 10}$  + Multiscale

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



10k points in 30-50ms



100k points in 100-200ms



## Geometric Loss functions for PyTorch

Our website: [www.kernel-operations.io/geomloss](http://www.kernel-operations.io/geomloss)

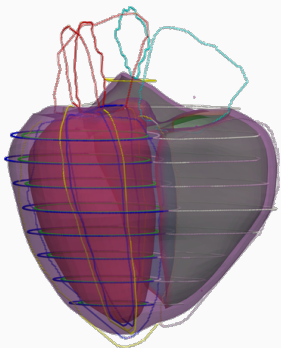
⇒ pip install geomloss ⇐

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

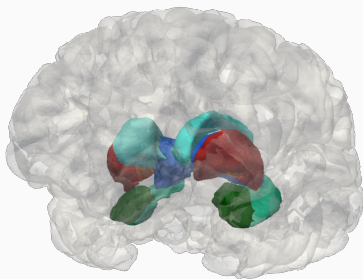
# Define a Wasserstein loss between sampled measures
from geomloss import SamplesLoss
loss = SamplesLoss(loss="sinkhorn", p=2)
L = loss(x, y) # By default, use constant weights
```

Soon: efficient support for **images**, **meshes** and generic metrics.

## My first motivation: computational anatomy



**Fast OT-based registration** with  
S. Joutard, X. Hao, A. Young from KCL,  
Z. Shen, M. Niethammer from UNC.



**Diffeomorphic and spline registration**  
e.g. Deformetrica LDDMM software  
with the Aramis Inria team.

KeOps and GeomLoss are:

- + **Fast:**  $\times 10$ - $\times 1,000$  speedup vs. naive GPU implementations.
- + **Memory-efficient:**  $O(N)$ , not  $O(N^2)$ .
- + **Versatile, with a transparent interface:** freedom!
- + **Powerful and well-documented:** research-friendly.
- Slow with **large vectors** of dimension  $D > 100$ .

## Future improvements

KeOps and GeomLoss are:

- + **Fast:**  $\times 10$ - $\times 1,000$  speedup vs. naive GPU implementations.
- + **Memory-efficient:**  $O(N)$ , not  $O(N^2)$ .
- + **Versatile, with a transparent interface:** freedom!
- + **Powerful and well-documented:** research-friendly.
- Slow with **large vectors** of dimension  $D > 100$ .

First half of 2021:

- **Approximation strategies** (Nyström, etc.) in KeOps.
- Wasserstein **barycenters** and **grid images** in GeomLoss.

### Roadmap for KeOps + GeomLoss:

- 2017–18 **Proof of concept** with conference papers, online codes.  
Get first feedback from the community.
- 2019–20 **Stable library** with solid theorems, a well-documented API.  
KeOps backends for high-level packages.
- 2021–22 **Mature library** with focused application papers, full tutorials.  
Works out-of-the-box for students and engineers.
- 2022+ **A standard toolbox**, with genuine clinical applications?  
That's the target!

# Conclusion

---

## Key points

- **Symbolic** matrices are to **geometric** ML what **sparse** matrices are to **graph** processing:
  - KeOps, **x30 speed-up** vs. PyTorch, TF and JAX.
  - Useful in a wide range of settings.
- Optimal Transport = **generalized sorting**:
  - Geometric gradients.
  - Super-fast  $O(N \log N)$  solvers.
- These tools open **new paths** for geometers and statisticians:
  - GPUs are more **versatile** than you think.
  - Ongoing work to provide **fast GPU backends** to researchers
  - going beyond what Google and Facebook are ready to pay for.

We believe that **KeOps** and **GeomLoss** will stimulate research on:

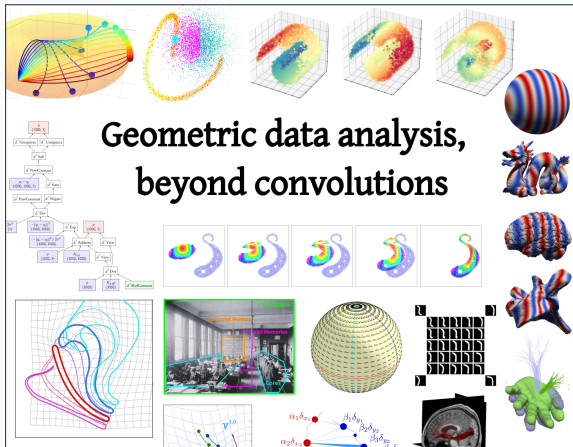
- **Clustering** methods: fast K-Means and EM iterations.
- Data **representation**: UMAP, fast KNN graphs with any metric.
- **Kernel** methods: kernel matrices.
- **Gaussian** processes: covariance matrices.
- **Geometric** deep learning: point convolutions.
- Geometric **statistics**: going beyond Euclidean models.
- Natural **language** processing: transformer networks?

What do you think?



# Documentation and tutorials are available online

⇒ [www.kernel-operations.io](http://www.kernel-operations.io) ⇐





Dimitri P Bertsekas.

**A distributed algorithm for the assignment problem.**

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



Grégoire Clarté, Antoine Diez, and Jean Feydy.

**Collective proposal distributions for nonlinear MCMC samplers: Mean-field theory and fast implementation.**

*arXiv preprint arXiv:1909.08988, 2019.*



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.*



Marco Cuturi.

**Sinkhorn distances: Lightspeed computation of optimal transport.**

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



Pierre Degond, Amic Frouvelle, Sara Merino-Aceituno, and Ariane Trescases.

**Alignment of self-propelled rigid bodies: from particle systems to macroscopic equations.**

In *International workshop on Stochastic Dynamics out of Equilibrium*, pages 28–66. Springer, 2017.



Pierre Degond and Sébastien Motsch.

**Continuum limit of self-driven particles with orientation interaction.**

*Mathematical Models and Methods in Applied Sciences*,  
18(supp01):1193–1215, 2008.



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.



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.



Florent Leclercq.

**Bayesian optimization for likelihood-free cosmological inference.**

*Physical Review D*, 98(6):063511, 2018.



Bruno Lévy.

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

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



Quentin Mérigot.

**A multiscale approach to optimal transport.**

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



Bernhard Schmitzer.

**Stabilized sparse scaling algorithms for entropy regularized transport problems.**

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



Freyr Sverrisson, Jean Feydy, Bruno E. Correia, and Michael M. Bronstein.

**Fast end-to-end learning on protein surfaces.**

*bioRxiv*, 2020.