# High performance computing in python

Sept 7, 2021: 11:00 - 12:30

## Leandro Parente

✉ leandro.parente@opengeohub.org
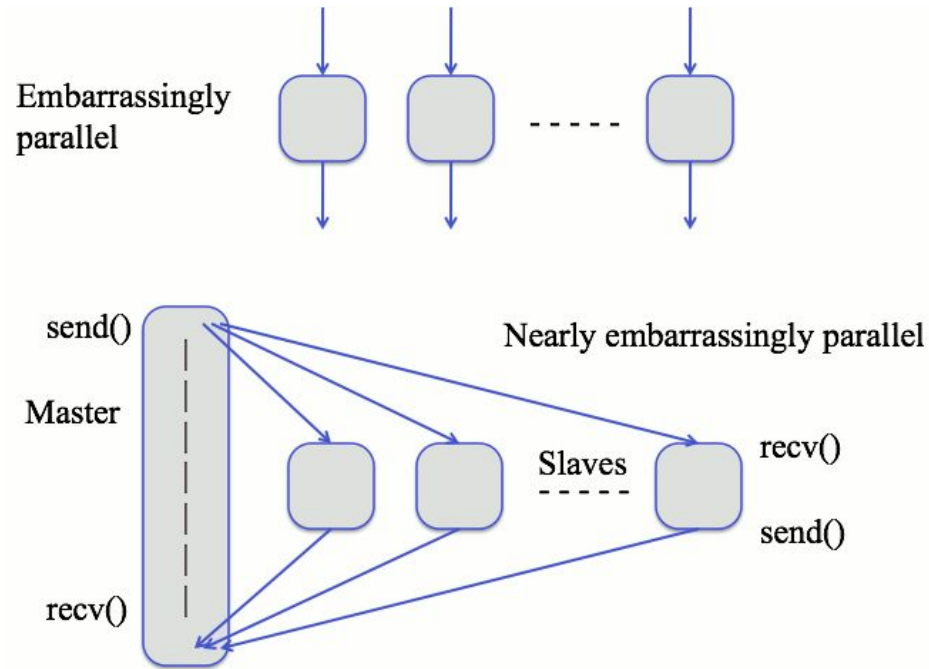
🏠 https://opengeohub.org

# Introduction to ODSE datasets in Python - Outline

- Embarrassingly parallel problems

- Possibilities to optimize a raster processing workflow

- BLAS and LAPACK implementations

- Optimizing a temporal array reduction and a numeric operations

- Production workflow using tile processing

OpenDataScience
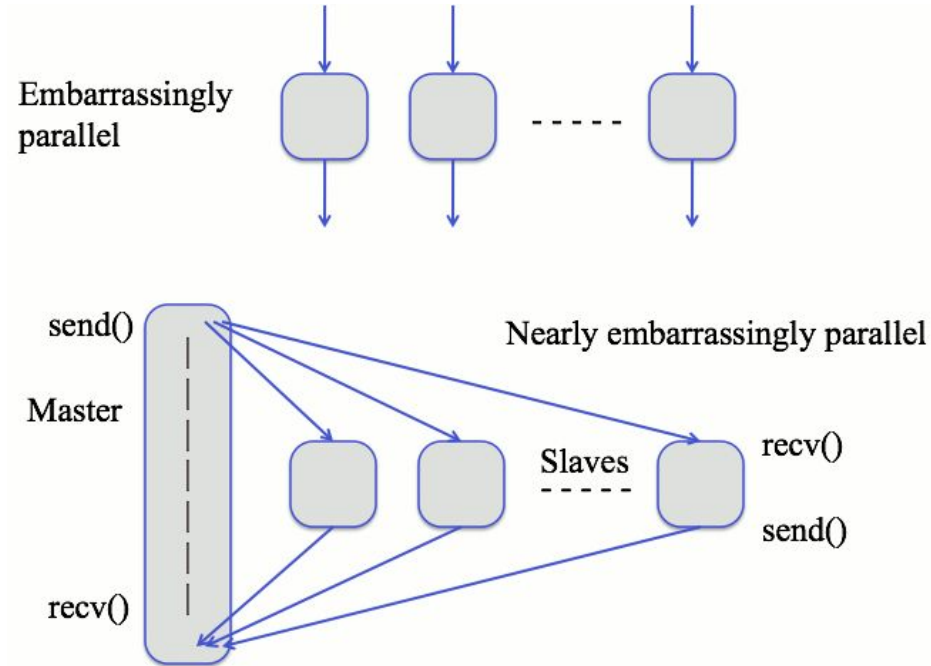
# Embarrassingly parallel problems

- Can be divided into completely **independent** parts,

- Requires none or very little communication,

- **Nearly embarrassingly parallel** is an embarrassingly parallel computation that requires initial data to be distributed and final results to be collected in some way
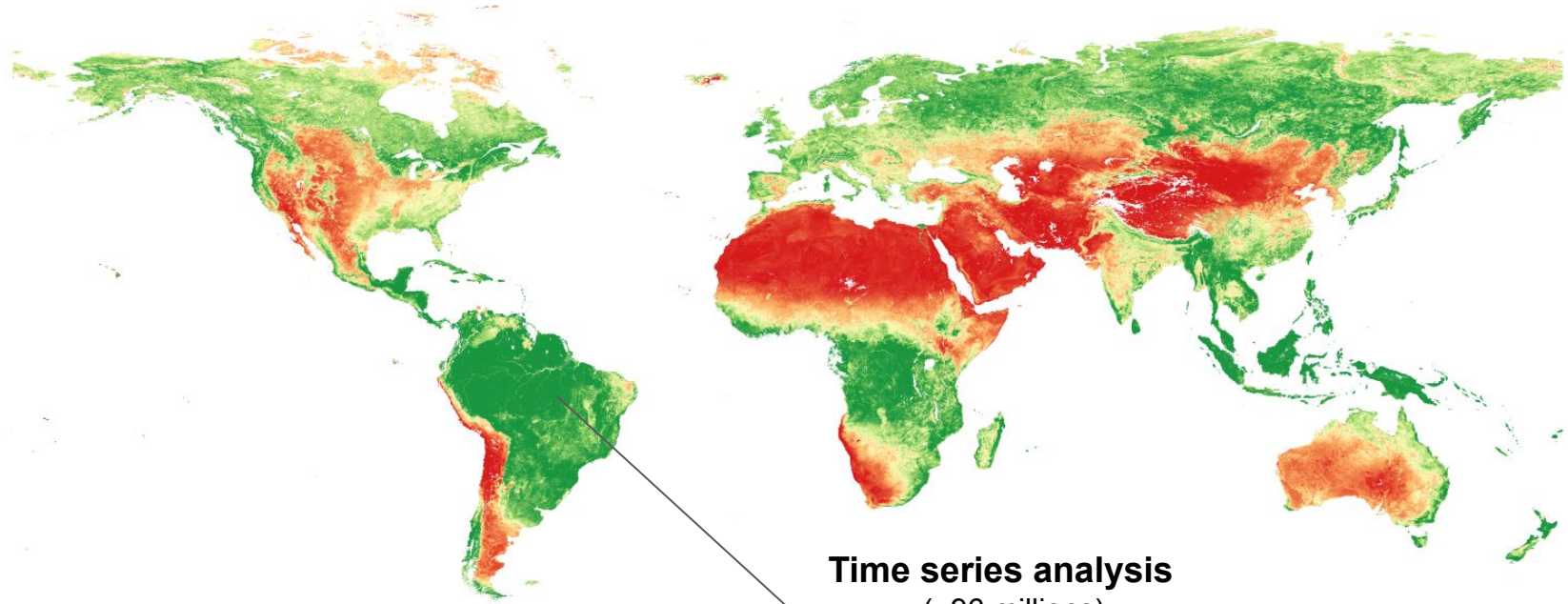


Source: &

**OpenDataScience**

# Embarrassingly parallel problems

- Can be divided into completely **independent** parts,

- Requires none or very little communication,

- **Nearly embarrassingly parallel** is an embarrassingly parallel computation that requires initial data to be distributed and final results to be collected in some way
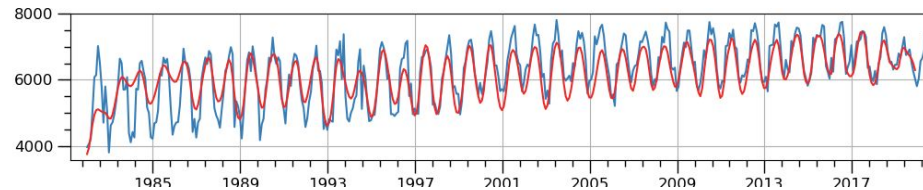
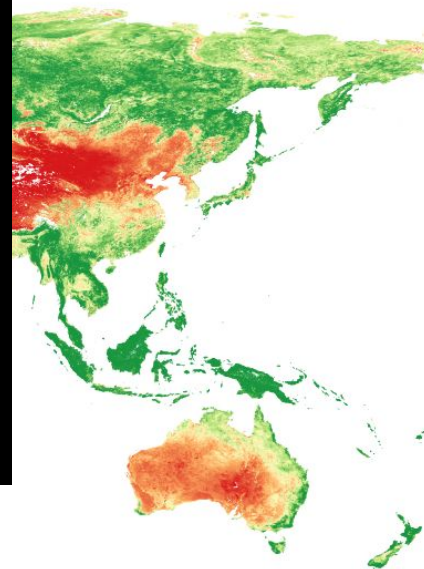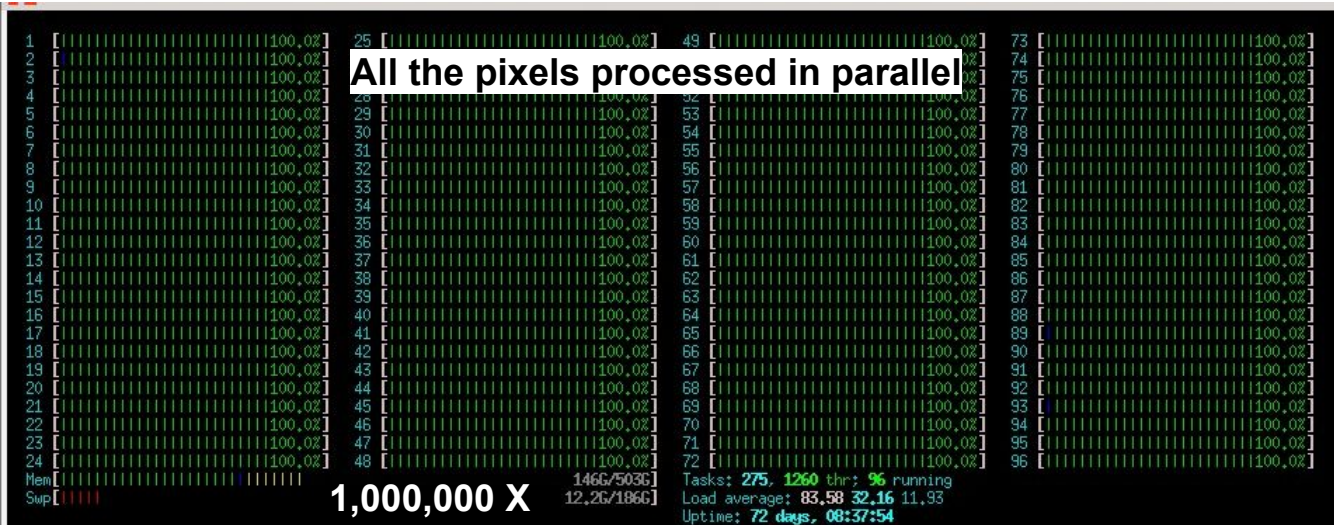**Raster processing**



Source: Embarrassingly Parallel Computations & Embarrassingly Parallel Algorithms

# Embarrassingly parallel problems



Time series analysis
(~96 millions)
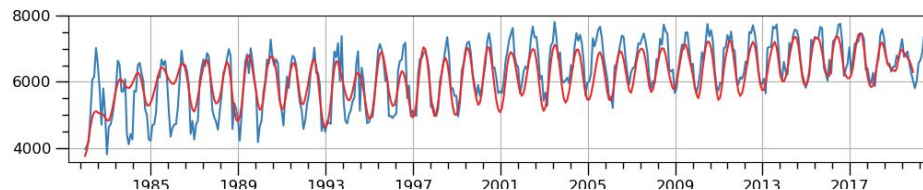
# Embarrassingly parallel problems



All the pixels processed in parallel

1,000,000 X

Time series analysis
(~96 millions)

# Possibilities to optimize a raster processing workflow

- Increase the number of CPU cores

- Improve data transfer speed

- Improve the processing code (new algorithms/functions):

# Possibilities to optimize a raster processing workflow

- Increase the number of CPU cores

- Improve data transfer speed

- Improve the processing code (new algorithms/functions):

  - Drop-in replacement

  - New code implementation
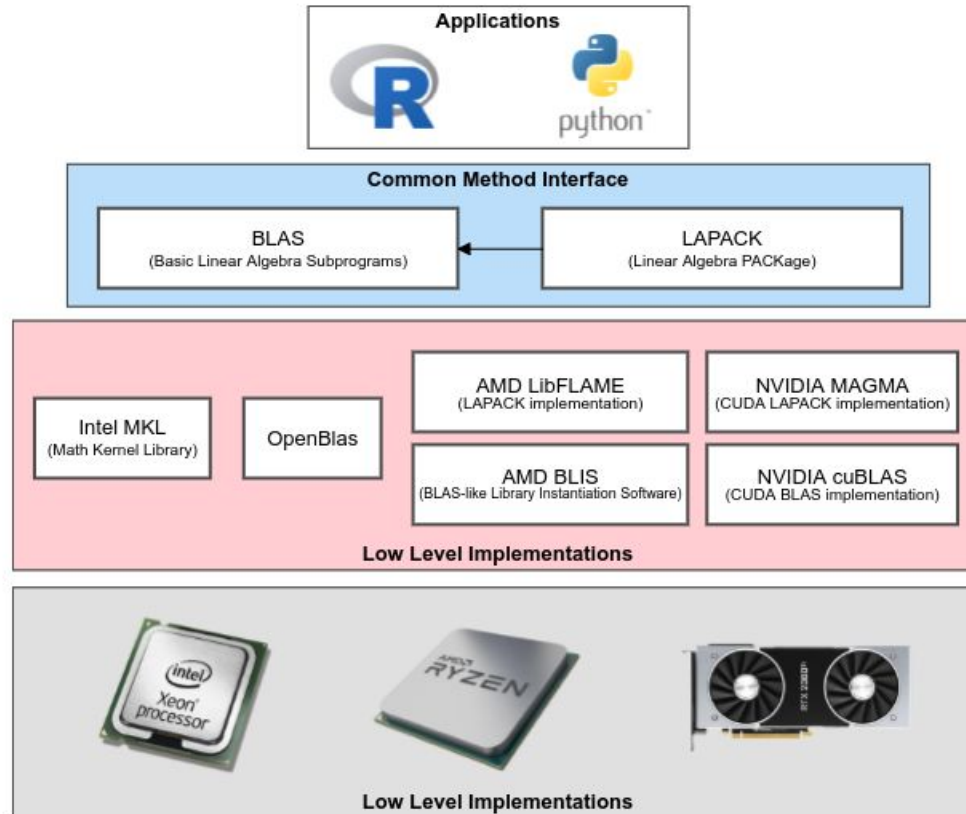
# BLAS and LAPACK implementations

BLAS (Basic Linear Algebra Subprograms) is a C library to provide a set of routines for basic vector and matrix operations

LAPACK (Linear Algebra PAckage) Fortran 90 library to solve linear equations, least-squares solutions of linear systems of equations, eigenvalue problems, singular value problems and the associated matrix factorization
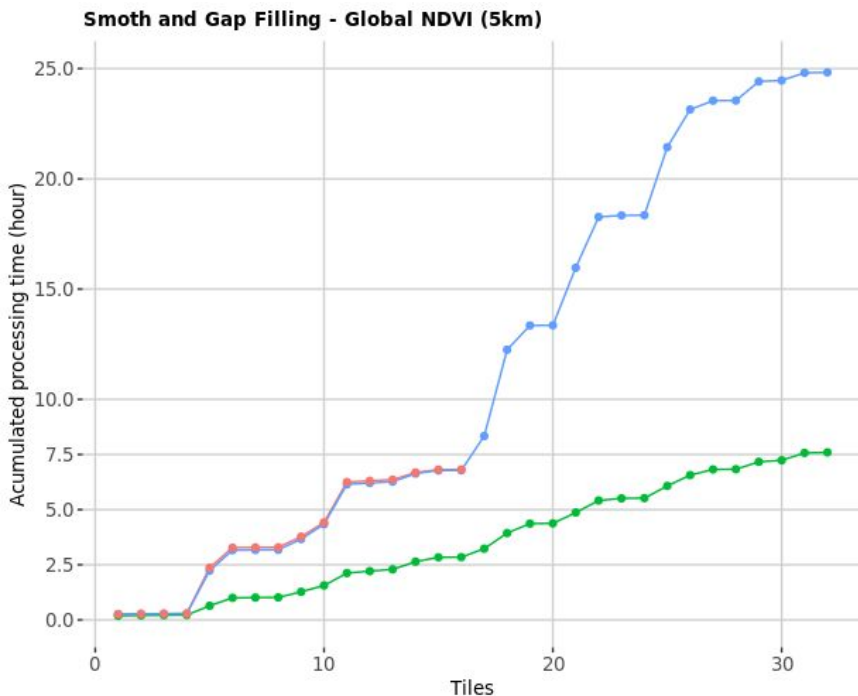


| Level 1 BLAS | | | | | | | prefixes |
|---|---|---|---|---|---|---|---|
| | dim scalar vector | vector | scalars | | 5-element array | | |
| SUBROUTINE xROTG ( | | | A, B, C, S ) | | | Generate plane rotation | S, D |
| SUBROUTINE xROTMG( | | | D1, D2, A, B, | | PARAM ) | Generate modified plane rotation | S, D |
| SUBROUTINE xROT  ( N, | X, INCX, Y, INCY, | | | C, S ) | | Apply plane rotation | S, D |
| SUBROUTINE xROTM ( N, | X, INCX, Y, INCY, | | | | PARAM ) | Apply modified plane rotation | S, D |
| SUBROUTINE xSWAP ( N, | X, INCX, Y, INCY ) | | | | | $x \leftrightarrow y$ | S, D, C, Z |
| SUBROUTINE xSCAL ( N, | ALPHA, X, INCX ) | | | | | $x \leftarrow \alpha x$ | S, D, C, Z, CS, ZD |
| SUBROUTINE xCOPY ( N, | X, INCX, Y, INCY ) | | | | | $y \leftarrow x$ | S, D, C, Z |
| SUBROUTINE xAXPY ( N, | ALPHA, X, INCX, Y, INCY ) | | | | | $y \leftarrow \alpha x + y$ | S, D, C, Z |
| FUNCTION    xDOT  ( N, | X, INCX, Y, INCY ) | | | | | $dot \leftarrow x^T y$ | S, D, DS |
| FUNCTION    xDOTU ( N, | X, INCX, Y, INCY ) | | | | | $dot \leftarrow x^T y$ | C, Z |
| FUNCTION    xDOTC ( N, | X, INCX, Y, INCY ) | | | | | $dot \leftarrow x^H y$ | C, Z |
| FUNCTION    xxDOT ( N, | X, INCX, Y, INCY ) | | | | | $dot \leftarrow \alpha + x^T y$ | SDS |
| FUNCTION    xNRM2 ( N, | X, INCX ) | | | | | $nrm2 \leftarrow \|x\|_2$ | S, D, SC, DZ |
| FUNCTION    xASUM ( N, | X, INCX ) | | | | | $asum \leftarrow \|re(x)\|_1 + \|im(x)\|_1$ | S, D, SC, DZ |
| FUNCTION    IxAMAX( N, | X, INCX ) | | | | | $amax \leftarrow 1^{st} k \ni \|re(x_k)\| + \|im(x_k)\|$ $= max(\|re(x_i)\| + \|im(x_i)\|)$ | S, D, C, Z |

# BLAS and LAPACK implementations

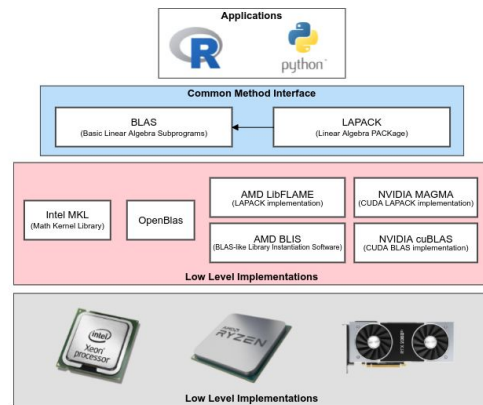# BLAS and LAPACK implementations

MKL is 3x faster then OpenBlas
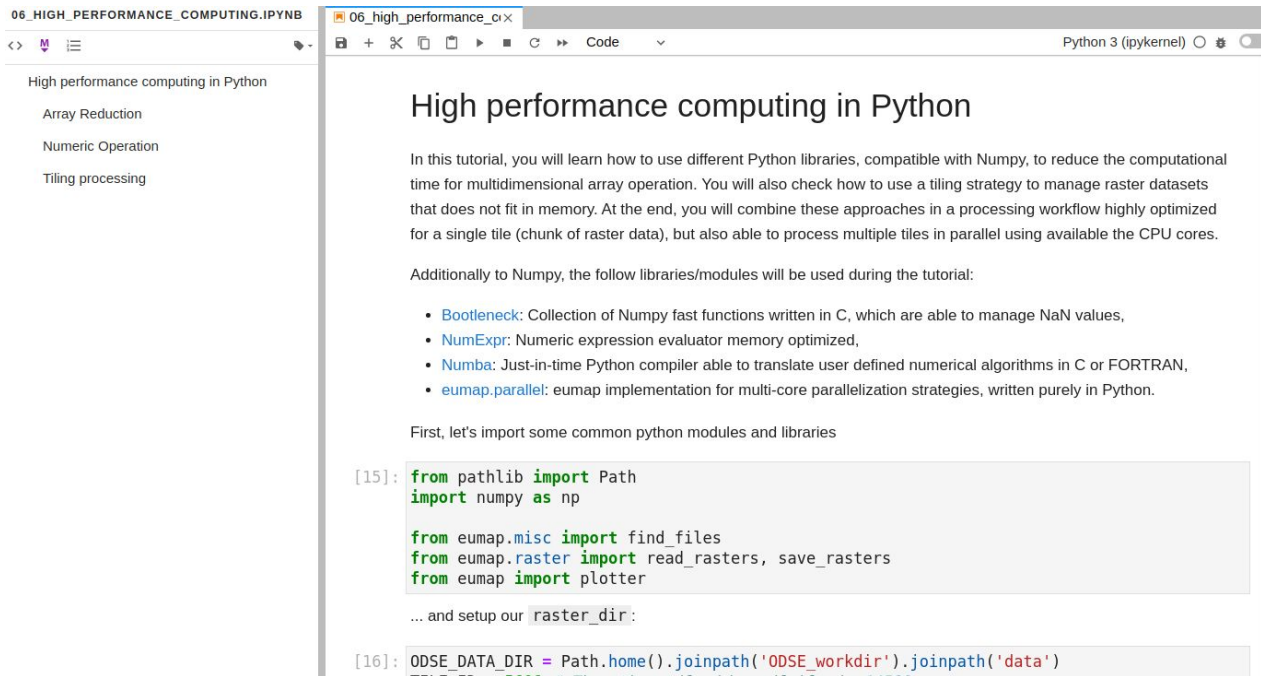
**Time series analysis**
(~96 millions)





Smoth and Gap Filling - Global NDVI (5km)

cat
- Blis+LibFLAME ((AMD EPYC 7702P - 108 Thr∈
- MKL (Intel Xeon Gold 6248 - 80 Threads)
- OpenBlas (AMD EPYC 7702P - 108 Threads)

OpenDataScience

# Optimizing a temporal array reduction and a numeric operations

https://gitlab.com/geoharmonizer_inea/odse-workshop-2021

# Production workflow using tile processing

# Production workflow using tile processing



1760 tiles

**CPU server 1**
(80 threads and 377 GB)

1760 tiles

**CPU server 2**
(96 threads and 504 GB)

1761 tiles

**CPU server 3**
(96 threads and 503 GB)

**Aggregate
the results**

1761 tiles

**CPU server 4**
(96 threads and 503 GB)

Data access

**SSD NAS**
(20 TB)

MINIO

S3

*Cloud optimized*

*tiff generation*

gdalbuildvrt
gdal_translate