

Data processing with XDS *demo - tutorial*

Kay Diederichs

kay.diederichs@uni-konstanz.de



CCP4@DLS 2024-11-25

Outline

- 1) Browsing through this presentation: Introduction to XDS (XSCALE, XDSCC12)
- 2) Demo: Processing data with XDSGUI (and using XDSSTAT, XDSCC12, SPOT2PDB, COOT, POINTLESS) of example data
- 3) Until Wed morning: process YOUR existing experimental data

The *XDS* program suite

- Original author: Wolfgang Kabsch (Max-Planck-Institute Heidelberg)
- Since ~1986
- I joined 2007



The XDS+ programs

- **XDS**: the main program - indexing, integrating, scaling, statistics
- **XSCALE**: scale several XDS intensity data sets together; zero-dose extrapolation; statistics
- **XDSCONV**: convert to MTZ / SHELX /... format (AIMLESS and CTRUNCATE are not needed!)

Programs independent from the XDS distribution:

- **XDS-Viewer**: inspect diagnostic images written by XDS, or (single) data frames (open source). *adxv* or *dials.image_viewer* may be used instead.
- **XDSTAT**: additional statistics
- **XDSGUI**: graphical user interface for XDS, SHELX C/D/E, ARCIMBOLDO (open source)
- **XDSCC12**: (XDS) which frames are bad?
(XSCALE) which data sets to re-index and merge?
- **XSCALE_ISOCLUSTER** multi-data-set: visualize relations and cluster

Sources of information

- XDS main website: <https://xds.mr.mpg.de> ; complete, accurate, up-to-date documentation; download: Linux, Intel- and Silicon-Mac; for Windows use WSL.
- XDSwiki: <https://wiki.uni-konstanz.de/xds>
Installation; data sets; documentation; download; links to e.g. Matthew J. Whitley's excellent tutorial given at CSHL 2018
- CCP4 bulletin board
- SBGrid talk (May 2020) at <https://www.youtube.com/watch?v=3WU9NrILECo>
- XDSGUI paper (2023) <https://doi.org/10.1107/S1600576723007057>
- Making a difference in multi-data-set crystallography: simple and deterministic data-scaling/selection methods. Assmann, G.M., Wang, M., Diederichs, K. (2020) Acta Cryst D76, 636 (serial crystallography, XDSCC12)

Automatic processing with XDS

- beamline software (provides **XDS.INP**)
- scripts: **xia2** (CCP4), **autoPROC** (Globalphasing), **generate_XDS.INP** (XDSwiki), **fast_dp** (Diamond, APS, NSLSII, ...), *xdsme* (Soleil), *autoxds* (SSRL), *autoprocess* (CMCF), ...
- CCP4: *pointless*, *xdsconv* (type CCP4_I+F, or CCP4, or CCP4_I, or CCP4_F)
- SHELX: *shelxc* reads XDS_ASCII.HKL

Principle of XDS processing

- There is one JOB= line in **XDS.INP** which specifies a list of tasks:

JOB= XYCORR INIT COLSPOT IDXREF DEFPIX INTEGRATE CORRECT

- data reduction is divided into tasks in a **modular** way
- information storage/exchange/flow between tasks by data files which may be inspected/analyzed
- each task needs the result from the previous tasks
- fine-tuning of a task does *not* require previous tasks to be repeated
- each task writes its output file **<TASK>.LP**

The tasks are ...

- XYCORR : write positional correction files
(**X-CORRECTIONS.cbf**, **Y-CORRECTIONS.cbf**)
- INIT : find background pixels (defaults usually OK)
- COLSPOT: find reflection positions
- IDXREF : "index" reflections; user may supply/choose spacegroup
- XPLAN [not required] : strategy for data collection
- DEFPIX : mask shadows on detector (use *XDSGUI!*)
- INTEGRATE : evaluates intensities on all frames, writes **INTEGRATE.HKL** and **FRAME.cbf**
- CORRECT : **scales**, rejects outliers, statistics, writes scaled, unmerged **XDS_ASCII.HKL** (and other files)

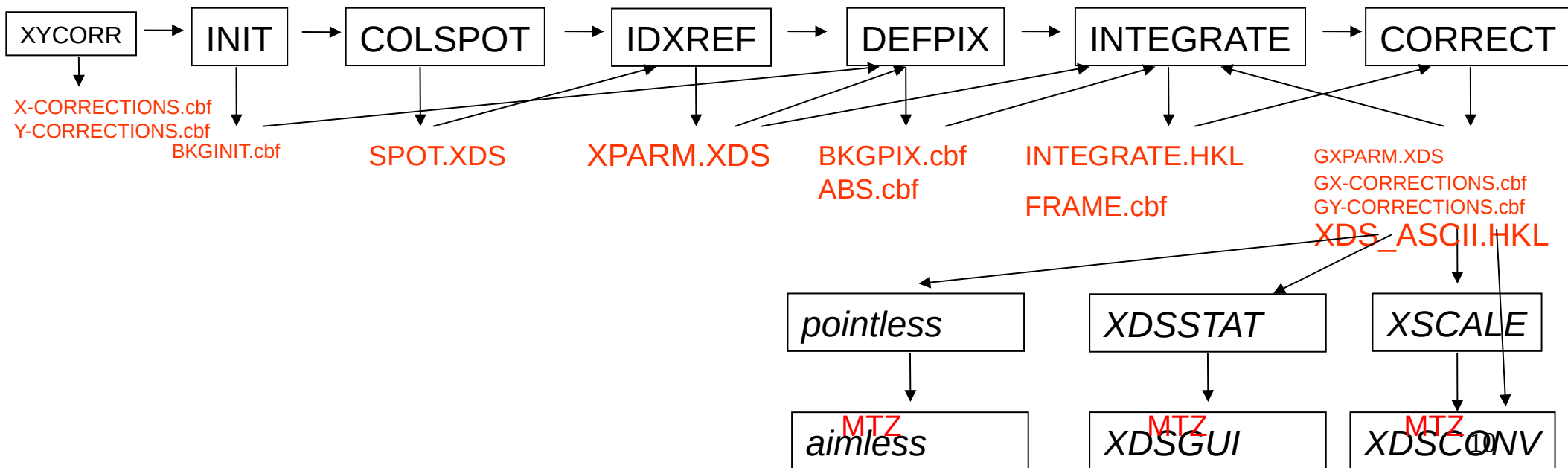
Example XDS.INP

```
JOB= XYCORR INIT COLSPOT IDXREF DEFPIX INTEGRATE CORRECT
ORGX=1546 ORGY=1552      !Detector origin (pixels); e.g. NX/2 NY/2
DETECTOR_DISTANCE=180    ! (mm)
OSCILLATION_RANGE=0.10   !degrees (>0)
X-RAY_WAVELENGTH=0.980243 !Angstroem
NAME_TEMPLATE_OF_DATA_FRAMES=frms/wga2-27_1_???.img
DATA_RANGE=1 3600        !Numbers of first and last data image collected
BACKGROUND_RANGE=1 10    !Numbers of first and last data image for background
SPACE_GROUP_NUMBER= 19   !0 for unknown crystals; cell constants are ignored.
UNIT_CELL_CONSTANTS= 44.4 86.4 104.5 90 90 90 ! not required if spgr=0
REFINE(IDXREF)=BEAM AXIS ORIENTATION CELL DISTANCE
REFINE(INTEGRATE)=DISTANCE BEAM ORIENTATION CELL ! AXIS
ROTATION_AXIS= 1.0 0.0 0.0
INCIDENT_BEAM_DIRECTION=0.0 0.0 1.0
FRACTION_OF_POLARIZATION=0.99                      ! SLS X06SA
POLARIZATION_PLANE_NORMAL= 0.0 1.0 0.0
DETECTOR=CCDCHESS      MINIMUM_VALID_PIXEL_VALUE=1      OVERLOAD=65000
DIRECTION_OF_DETECTOR_X-AXIS= 1.0 0.0 0.0
DIRECTION_OF_DETECTOR_Y-AXIS= 0.0 1.0 0.0
MINIMUM_NUMBER_OF_PIXELS_IN_A_SPOT=3 ! 3 for very sharp reflections, e.g. 6 for ordinary/bad data sets
VALUE_RANGE_FOR_TRUSTED_DETECTOR_PIXELS= 7000 30000 !Used by DEFPIX
                                     !for excluding shaded parts of the detector.
INCLUDE_RESOLUTION_RANGE=50.0 1.3 !Angstroem; used by DEFPIX,INTEGRATE,CORRECT
```

Bold keyword/parameter pairs are required; **yellow** ones change between experiments; **this** may need adjustment since it depends on the crystal. Documentation: xds.mr.mpg.de/html_doc/xds_parameters.html

Information flow

NAME_TEMPLA	OSCILLATION_	ORGX	
TE_OF_DATA_	RANGE	ORGY	
FRAMES	SEPMIN	DETECTOR_DISTANCE	DATA_RANGE
DETECTOR	STRONG_PIXEL	X_RAY_WAVELENGTH	
		SPACE_GROUP_NUMBER	



Example **XSCALE.INP**

```
!===== EXAMPLE 3: specific reindexing of input data sets
!
!      Use of specific reindexing of input data sets for resolving
!      indexing ambiguities in the scaled output data set. This
!      happens if the crystal's space group symmetry is lower than
!      its lattice symmetry.
!
RESOLUTION_SHELLS= 100 10 6 4 3 2 1.9
SPACE_GROUP_NUMBER=78
UNIT_CELL_CONSTANTS=57.39 57.39 106.9    90 90 90
OUTPUT_FILE=scaf8_all_merged.hkl
MERGE=TRUE FRIEDEL'S_LAW=FALSE
STRICT_ABSORPTION_CORRECTION=TRUE
INPUT_FILE= ../xds-1_2/XDS_ASCII.HKL
REIDX_ISET= -1  0  0  0  0  1  0  0  0  0 -1  0
INPUT_FILE= ../xds-2_1/XDS_ASCII.HKL
INPUT_FILE= ../xds-3_1/XDS_ASCII.HKL
INPUT_FILE= ../xds-1_4/XDS_ASCII.HKL
INPUT_FILE= *../xds-5_1/XDS_ASCII.HKL
```

Bold keyword/parameter pairs are required. Complete documentation at
xds.mr.mpg.de/html_doc/xscale_parameters.html

Output file **XSCALE.LP** shows level of systematic error **ISa** before/after scaling.

Example XDSCONV.INP

```
! UNIT_CELL_CONSTANTS= 10 20 30 90 90 90  
! SPACE_GROUP_NUMBER= 96  
! GENERATE_FRACTION_OF_TEST_REFLECTIONS=0.05
```

INPUT_FILE=XDS_ASCII.HKL

OUTPUT_FILE=temp.hkl CCP4_I+F ! or CCP4_I or CCP4_F or SHELX or CNS

FRIEDEL'S_LAW=FALSE ! store anomalous signal in output file even if weak

Bold keyword/parameter pairs are required. Complete documentation at
xds.mr.mpg.de/html_doc/xdsconv_parameters.html

The signal and the noise: random and systematic errors

Random error

True randomness occurs due to the quantum nature of matter.

- counting photons
- electronic noise (detector, electronics)

Random error is proportional to square root of measured value

Level of random error is given by R-values, $CC_{1/2}$ and I/σ (Data quality assessment lecture)

Systematic error

- crystal: conditions, composition, conformation, damage due to experiment, ...
- apparatus: shadows, absorption, vibrations, photon/electron flux ...
- processing software: inaccurate or incomplete modelling of experiment

Systematic error is proportional to measured value (often 1..10% but sometimes much more e.g. in case of shadows and overloads)

Level of systematic error is given by $ISa = \text{asymptotic signal/noise}$ (CORRECT.LP, AIMLESS, dials.scale)

For important data, do not rely on automatic data processing

Synchrotrons typically run multiple pipelines, and the user has to choose ...

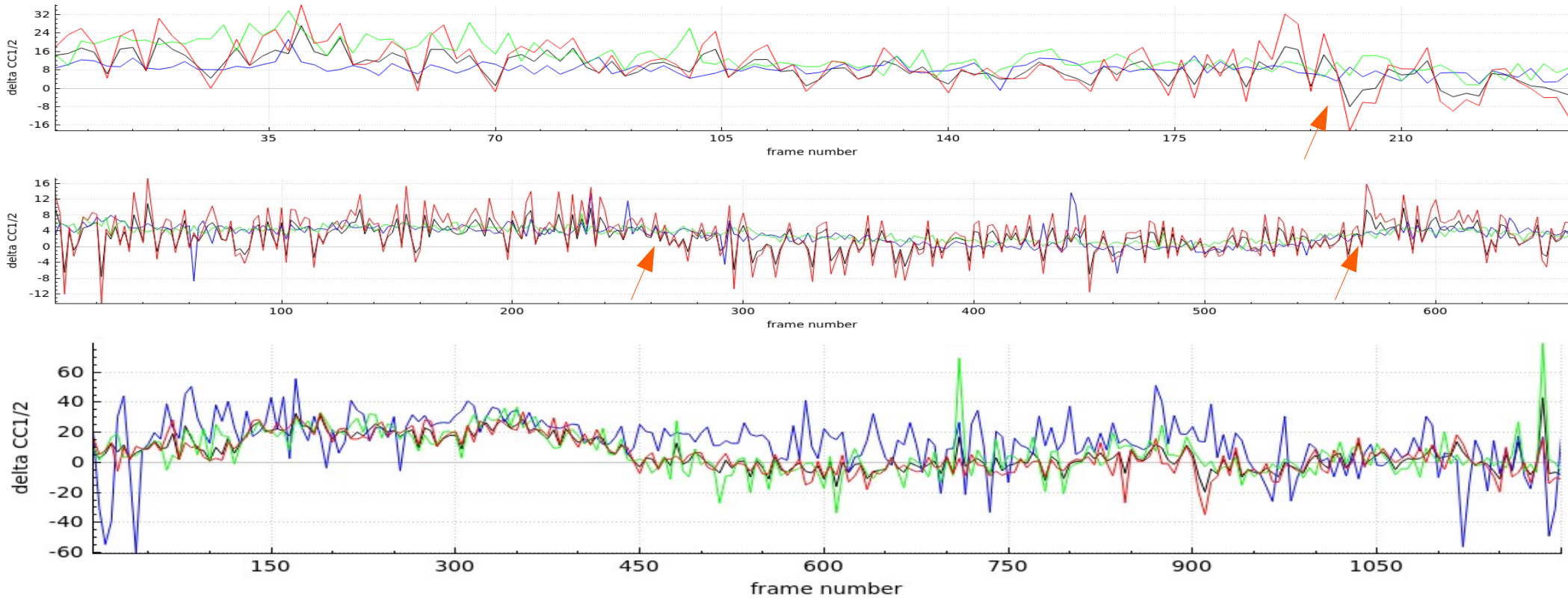
What can go wrong in automatic data processing?

- does not handle radiation damage (we need to discard frames towards the end of the data set, but where should we cut?)
- does not handle shadowed areas of detector
- does not handle indexing problems (multiple lattices, ice, ...) flexibly
- does not optimize processing
- there are no perfect criteria for the quality of a data set - “Table 1” does not tell the whole story

Automatic processing with GlobalPhasing's *autoPROC*, CCP4's *xia2*, ... is rather reliable for good data!

Difficult data sets typically benefit from human insight.

3 examples of single data sets (plots from *XDSGUI*)



XDSCC12: calculates $\Delta CC_{1/2,i} = CC_{1/2,with_i} - CC_{1/2,without_i}$

Three resolution ranges (blue=low green=medium red=high) - i refers to batches of width 1°

- find bad frame ranges
- radiation damage

Manual processing with XDSGUI

- problems in phasing and refinement: often due to bad / wrong data processing – reprocess raw data!
- visually inspect frames; mask shadows
- optimize processing parameters, frame range, resolution cutoff ..
- understand experiment: presentation of tables as plots; **tools** tab/Further analyses/**show spots in reciprocal space**
- user-extensible / modifiable commands

Optimize data processing

- XDSwiki:Optimisation#Re-INTEGRATEing_with_the_correct_spacegroup.2C_refined_geometry_and_fine-slicing_of_profiles

XDSwiki:Optimisation#using_the_refined_values_for_beam_divergence_and_d_mosaicity_for_re-integration

- **tools** tab:
 - * “Saving and comparing good results” and “Optimizing data quality”. After changing parameters, run “JOB=DEFPIX INTEGRATE CORRECT”, compare and save if better/restore old if worse.
 - * “Further analysis”: Inspect indexed/unindexed spots in reciprocal space
- consider use of *StarAniso* if anisotropy is strong (i.e. visible)

... may make the difference between structure solved or not, interpretable or non-interpretable map, good or bad refinement, ...

Summary

Data processing is the crucial link between experiment and structure

- * garbage in – garbage out!
- * for important data or in case of downstream difficulties, do not rely on automatic data processing
- * manual checks are easily performed with *XDSGUI*
- * try to optimize data processing – this converts noise to signal, and may enable structure solution, and/or improve refinement