# Serial crystallography
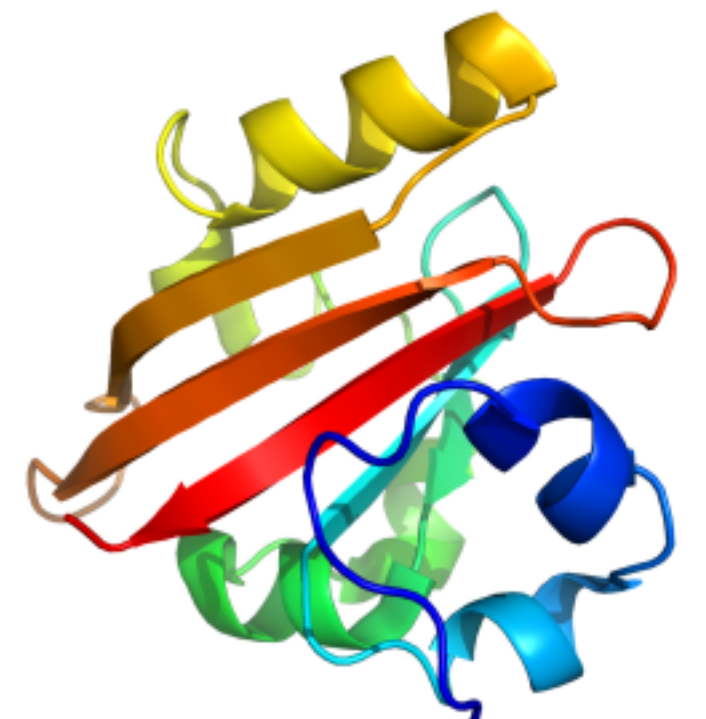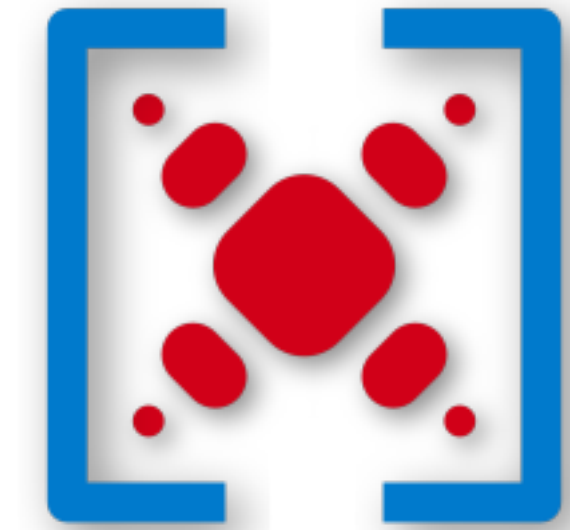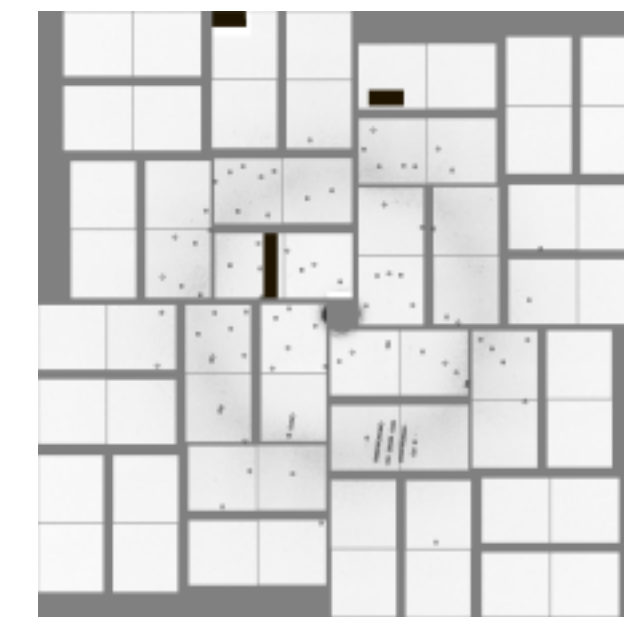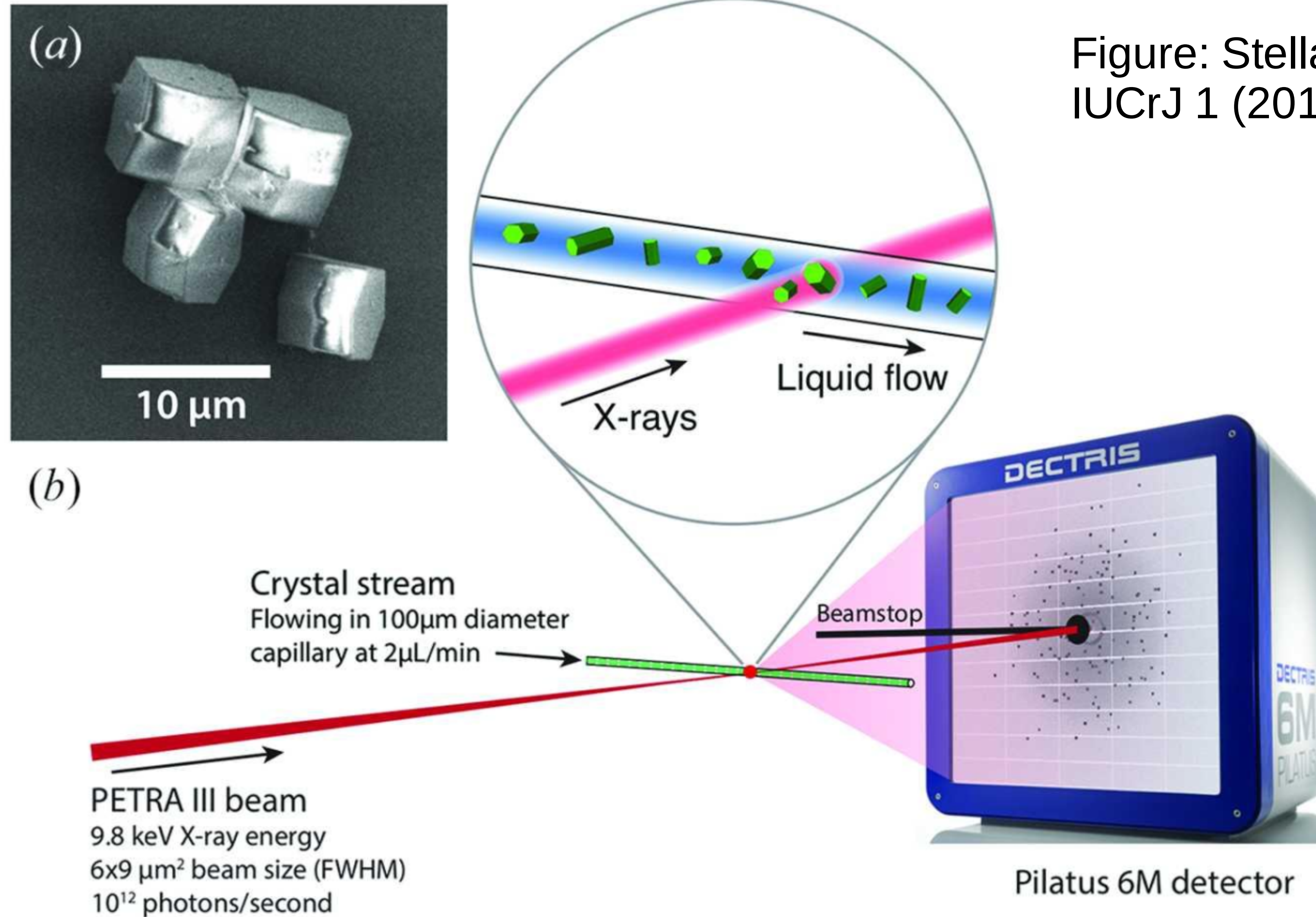
Data processing and thoughts about "The Digital
Scientific Method"

Thomas White
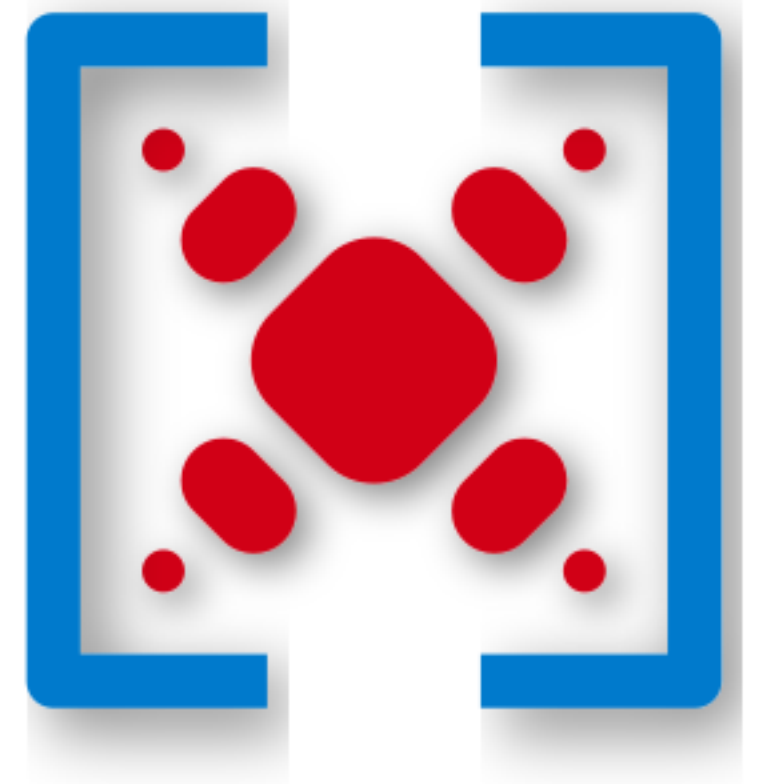7. MT Meeting
3. Feb 2021

# Serial crystallography



Figure: Stellato et al., IUCrJ 1 (2014) p204-212

(a)

10 μm

(b)

X-rays

Liquid flow

Crystal stream
Flowing in 100μm diameter
capillary at 2μL/min

Beamstop

PETRA III beam
9.8 keV X-ray energy
6x9 μm² beam size (FWHM)
$10^{12}$ photons/second

Pilatus 6M detector

# CrystFEL: data processing for serial crystallography



https://www.desy.de/~twhite/crystfel

Throughout 2019 and 2020, a new application of CrystFEL came out in peer-reviewed literature every 10 days (on average).
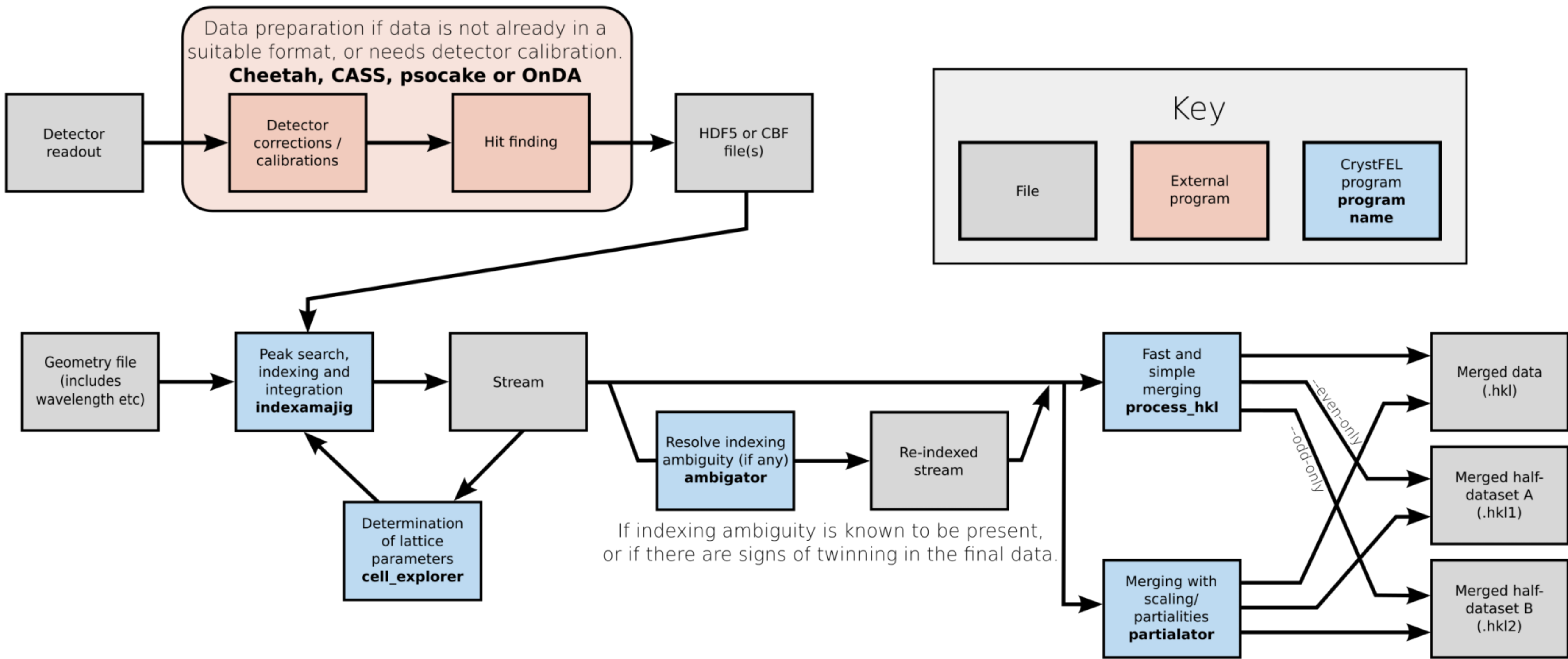
Development since 2009 at DESY.
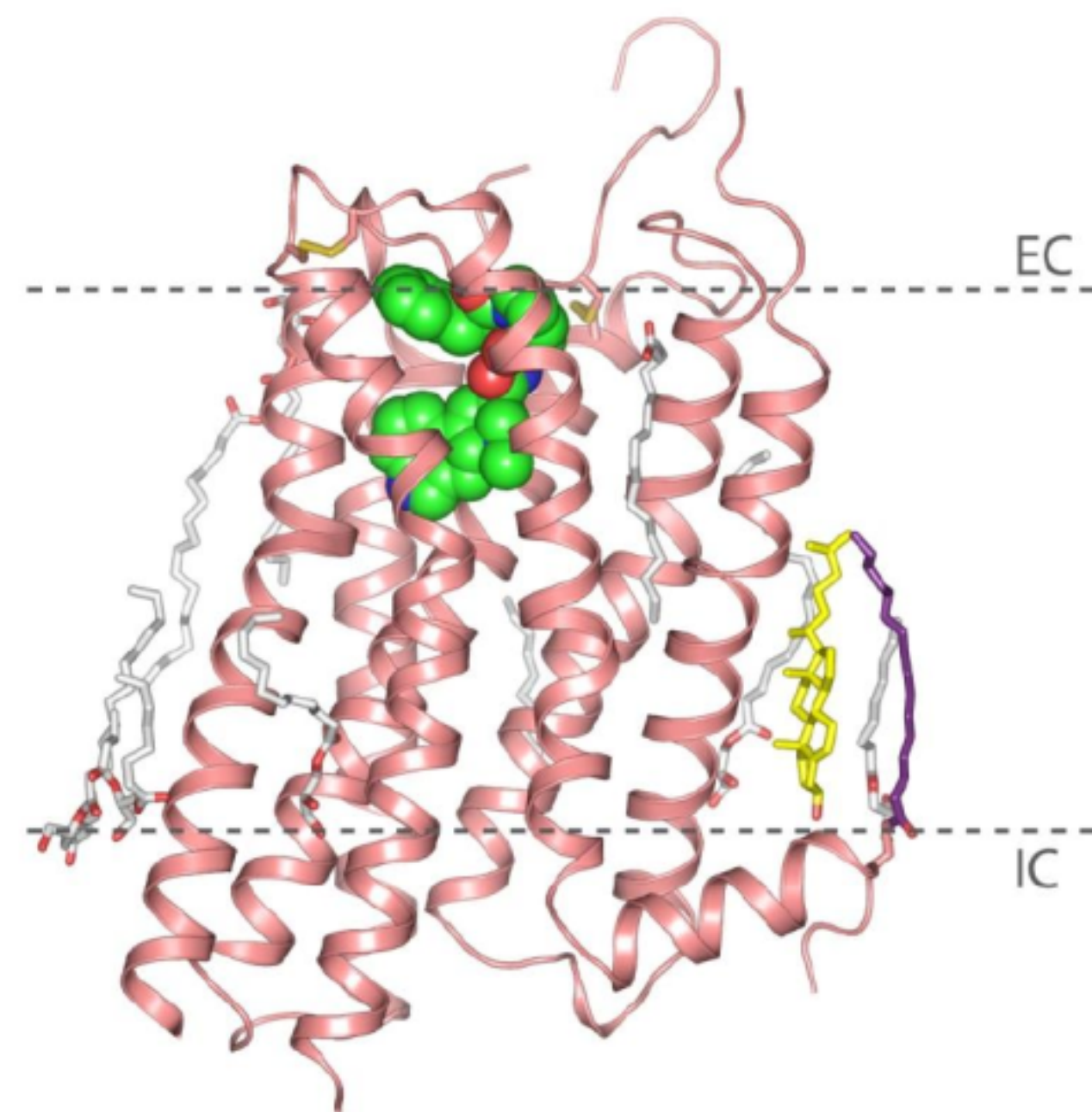
First "public" release in 2012.

Project led by me (Thomas White)

Tools for indexing, integration, merging, calculating figures of merit, simulating test data and more.

# CrystFEL processing pipeline

# What can we expect "users" to know?


Structural biology


Unix command line


Data formats


Crystallographic theory


Cluster / batch system


CrystFEL

# The CrystFEL GUI

## Peak search

Peak search algorithm: Local background estimation (peakfinder9)

Minimum SNR of brightest pixel in peak: 7.00

Minimum SNR of peak pixel: 6.00

Minimum signal/noise ratio: 5.00

Minimum background standard deviation: 11.00

Brightest pixel cutoff (just for speed): -inf

Local background radius: 3

Discard changes     Confirm

## Index all frames

Job name    Indexing    Integration    **Cluster/batch system**

Batch system: SLURM

Submit job to partition: maxwell

Split job into blocks of 1000 frames

Send notifications to: myself@example.org

Search path for executables: /path/to/indexing/pro...

Cancel     Run

---

CrystFEL

File   View   Tools   Help

...ents/5HT2B/r0092-ab/data1/LCLS_2013_Mar22_r0092_035038_16020.h5 // (frame 1 of 70)

Show results from: Calculations within GUI

Load data

Peak detection

Index this frame

Index all frames

Determine unit cell

Indexing ambiguity

Merge

Indexing all frames (5ht2b-index-all)

Cancel

Multi-lattice indexing ("delete and retry"): on

Retry indexing: on

WARNING: 53 reflections could not be integrated

WARNING: 49 reflections could not be integrated

WARNING: 30 reflections could not be integrated

WARNING: 26 reflections could not be integrated

WARNING: 31 reflections could not be integrated

# Benefits of real-time data processing

Faster results and publication

No need to store raw data (except for audit/accountability)

Better situational awareness during experiment ("Do we have enough data yet?")

Faster diagnosis of experiment problems (e.g. pump laser misaligned)

Less scope for "fiddling and taking the result you like best"

# Serial crystallography



Figure: Stellato et al., IUCrJ 1 (2014) p204-212

(a)

10 µm

(b)

X-rays

Liquid flow

Crystal stream
Flowing in 100µm diameter
capillary at 2µL/min

Beamstop

PETRA III beam
9.8 keV X-ray energy
6x9 µm² beam size (FWHM)
$10^{12}$ photons/second

Pilatus 6M detector

# Benefits of real-time data processing

Faster results and publication

No need to store raw data (except for audit/accountability)

Better situational awareness during experiment ("Do we have enough data yet?")

Faster diagnosis of experiment problems (e.g. pump laser misaligned)

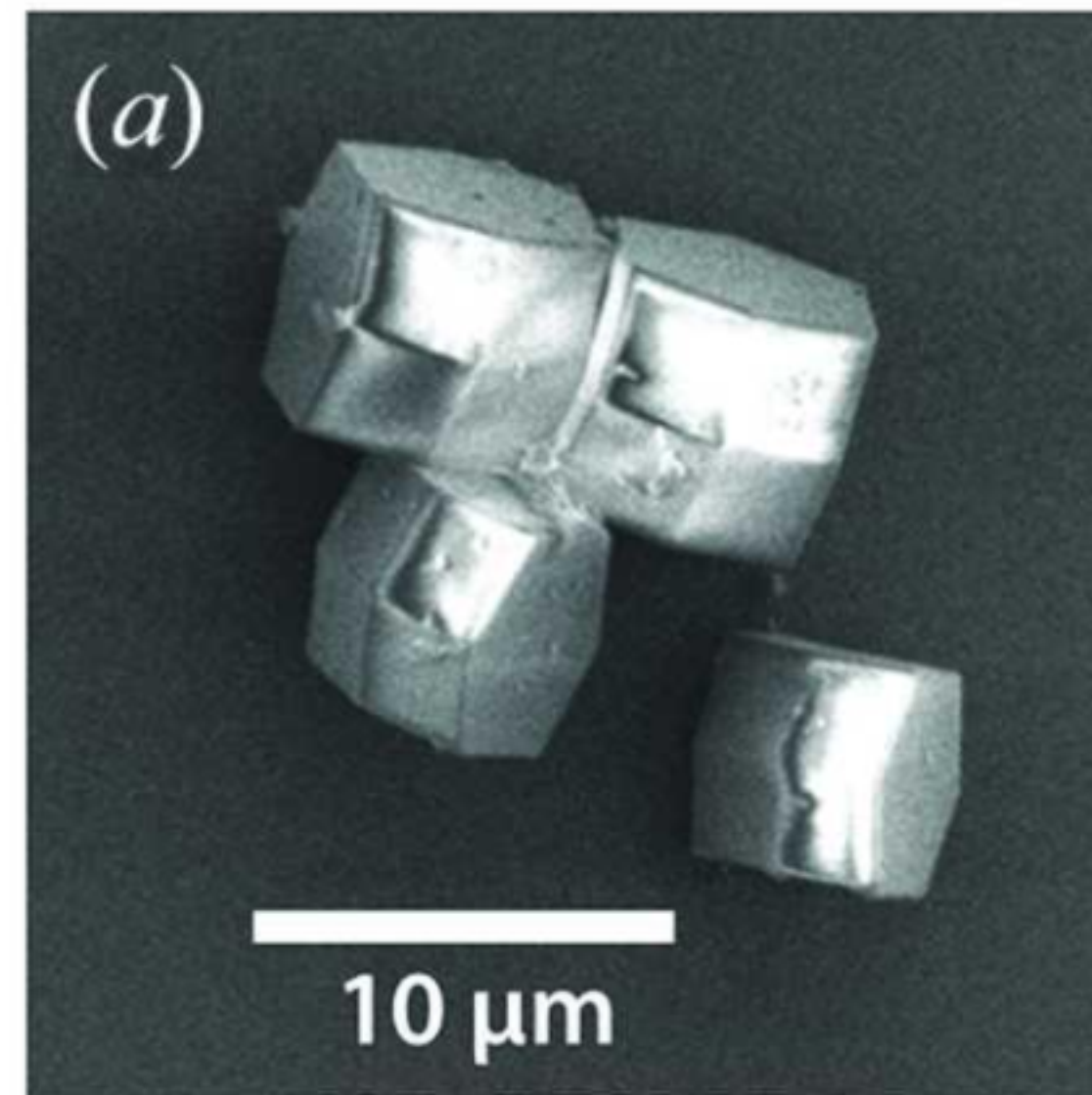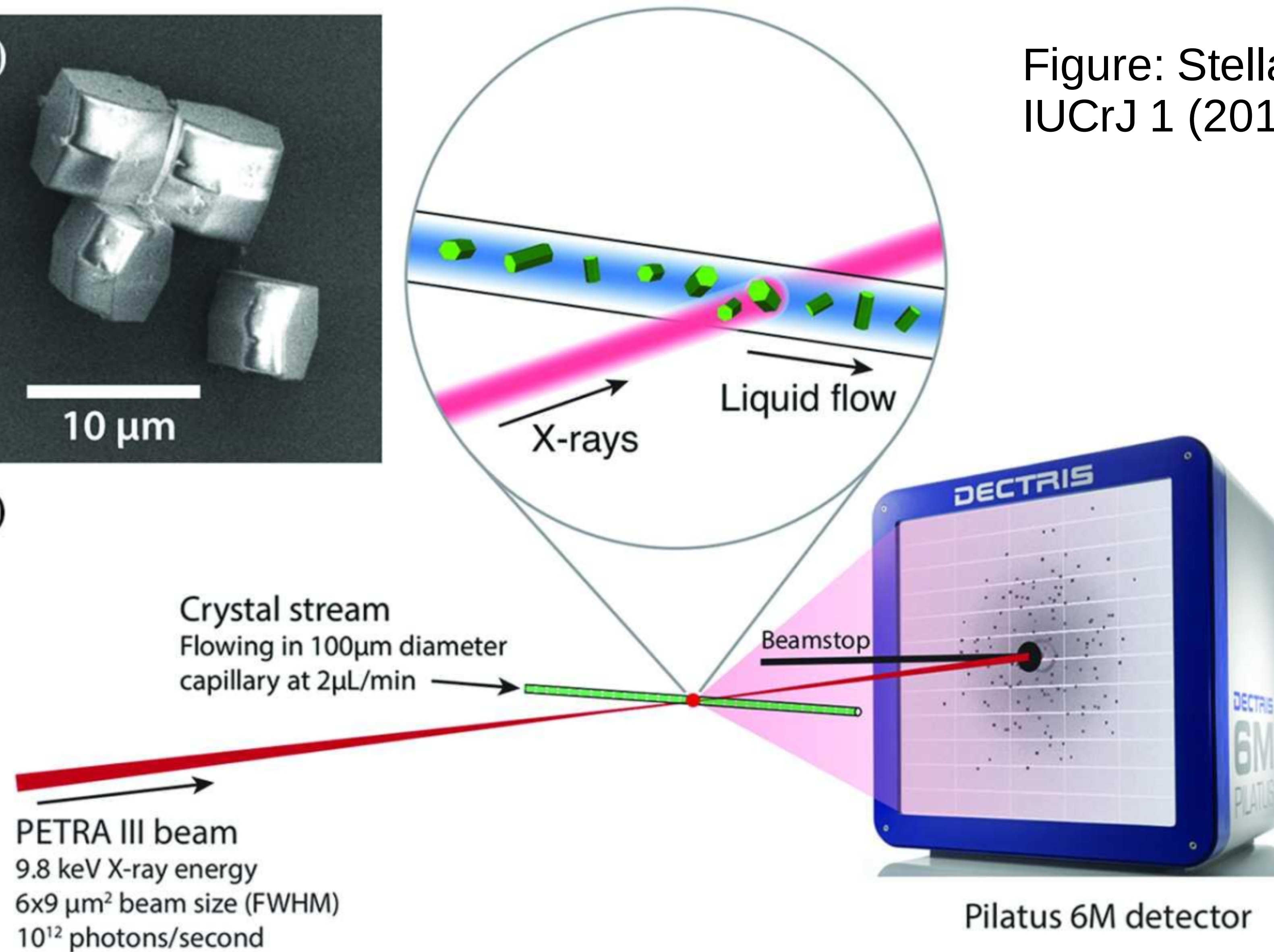Less scope for "fiddling and taking the result you like best" ——————>



Image: XKCD (2400)

# Towards real-time data processing

Remove bottlenecks, "make it go faster", e.g. by algorithm tuning:

# Further topics

Are there better ways to write code?

- Better level of abstraction
- Consider "old" computer science concepts,
    e.g. functional or declarative programming style

Making code part of the scientific record:

- Code vs. written methods section
- Notebook-oriented user interface
    (... for large batch processes?)
- Recognition of code as important scientific output
- Institutional software repository and DOIs for code

# Summary / take "home" points

To me, "The Digital Scientific Method" is about how we can use computers to improve the reliability and efficiency of our science.

Carefully considered user interfaces can make a big difference
- Embed domain-specific knowledge so that not everyone has to learn it
    (.... reduce the chance to make mistakes ....)
- Reach into the "guts" of the data processing software
- Close feedback loops in data processing workflow

Real-time processing will make a huge difference
- Speed
- Don't store raw data
- Better situational awareness during experiment
- Faster diagnosis of problems
- Less scope for selecting "nice" results   ——> more reliable science