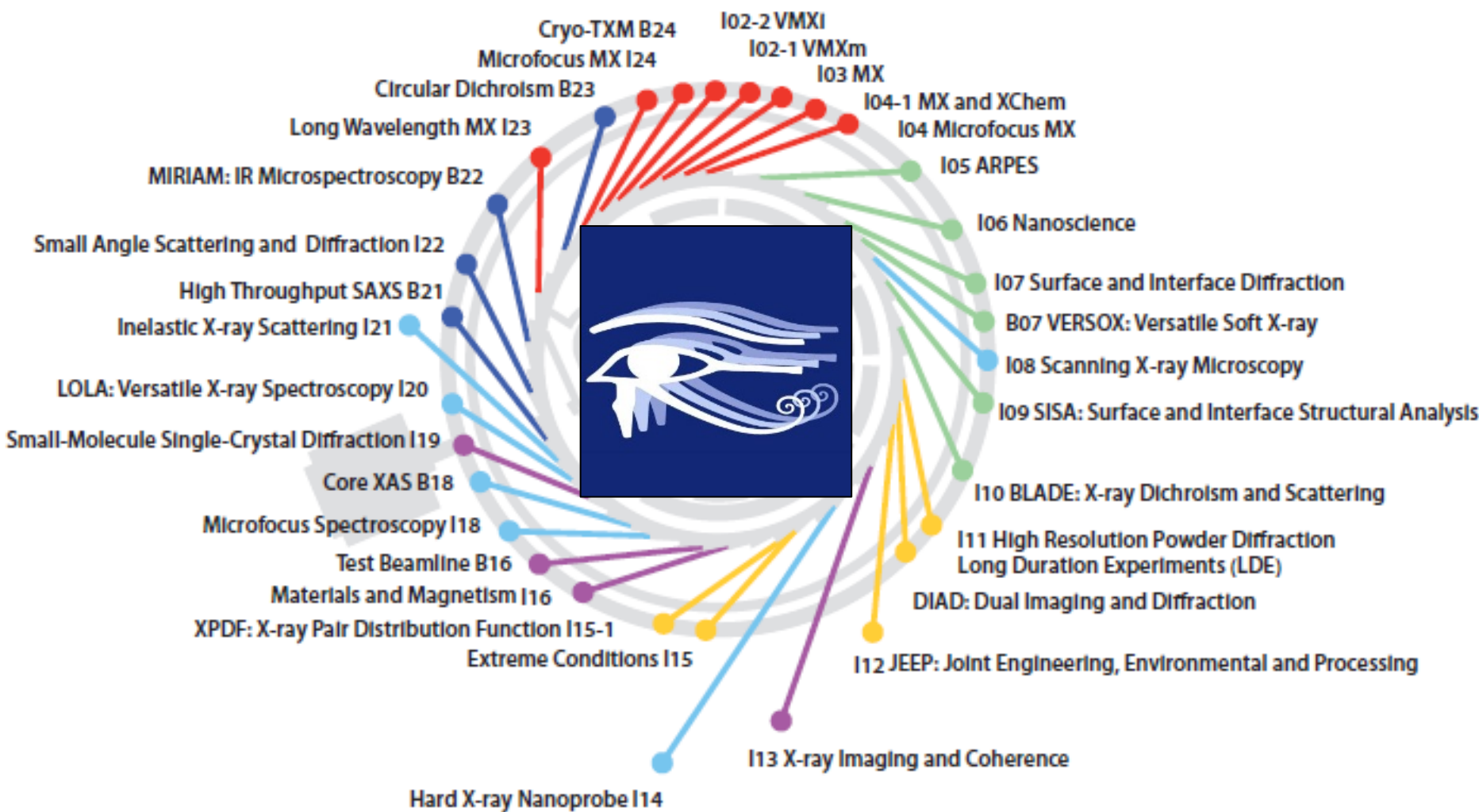


Plans following September Workshop



Neil A Smith
31st October 2019





Background

- ISPyB developer group holds (mostly) monthly Video conferences
- Changes to the database schema are initially proposed and captured on github
 - <https://github.com/ispyb/ispyb-database-modeling>
 - 31 Open, 16 Closed, housekeeping required
- Small changes to database schema are discussed and sometimes agreed
- Large changes (refactoring, quality improvements etc.) often deferred
- Agreed to hold technical workshop in September hosted by SOLEIL
 - Thank you to Tatiana Isabet and Idrissou Chado

Aim

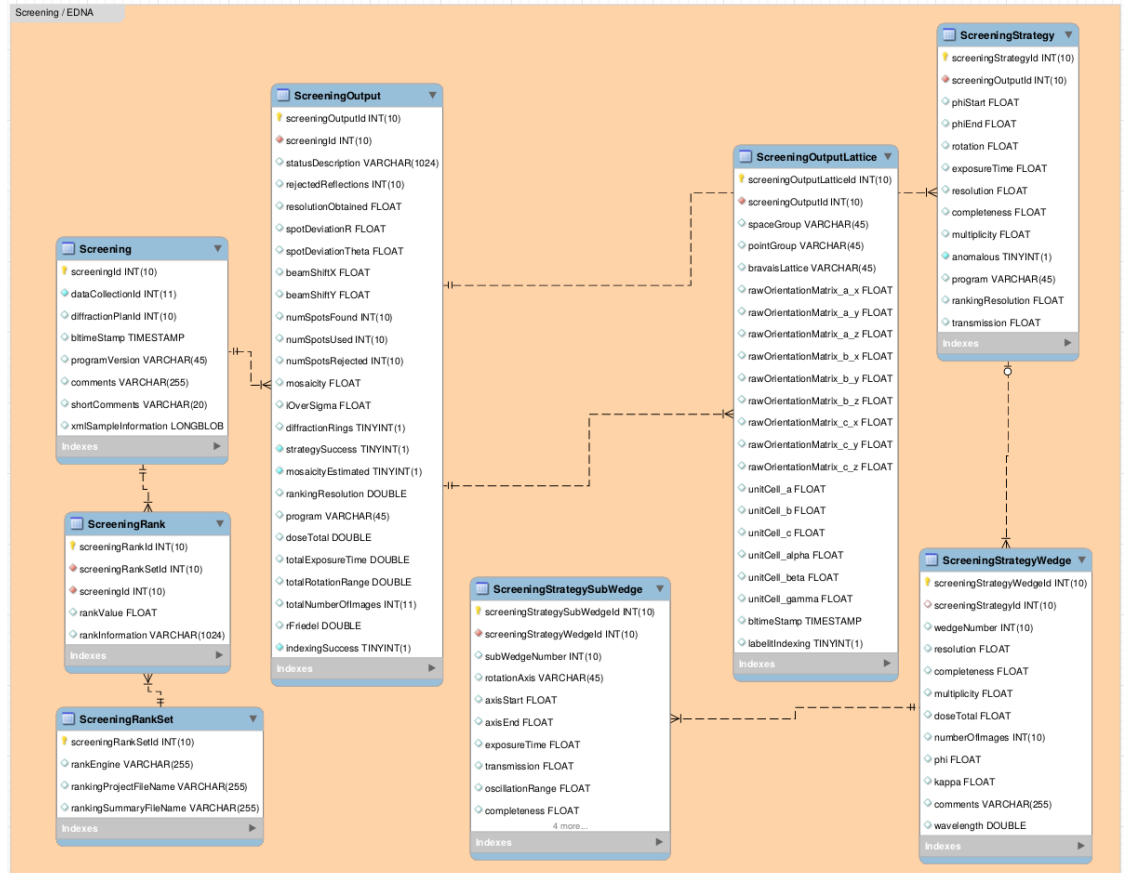
- Three database topics covered
 - Screening Tables
 - Data Collection
 - Sample tables
- Cover other items as time/energy allow!
- Details reported in meeting minutes on collaboration pages

https://ispyb.github.io/ISPyB/webpages/Other_meetings/TechMeeting_20190912_final.pdf

Thank you to Rasmus Fogh for compiling the minutes!

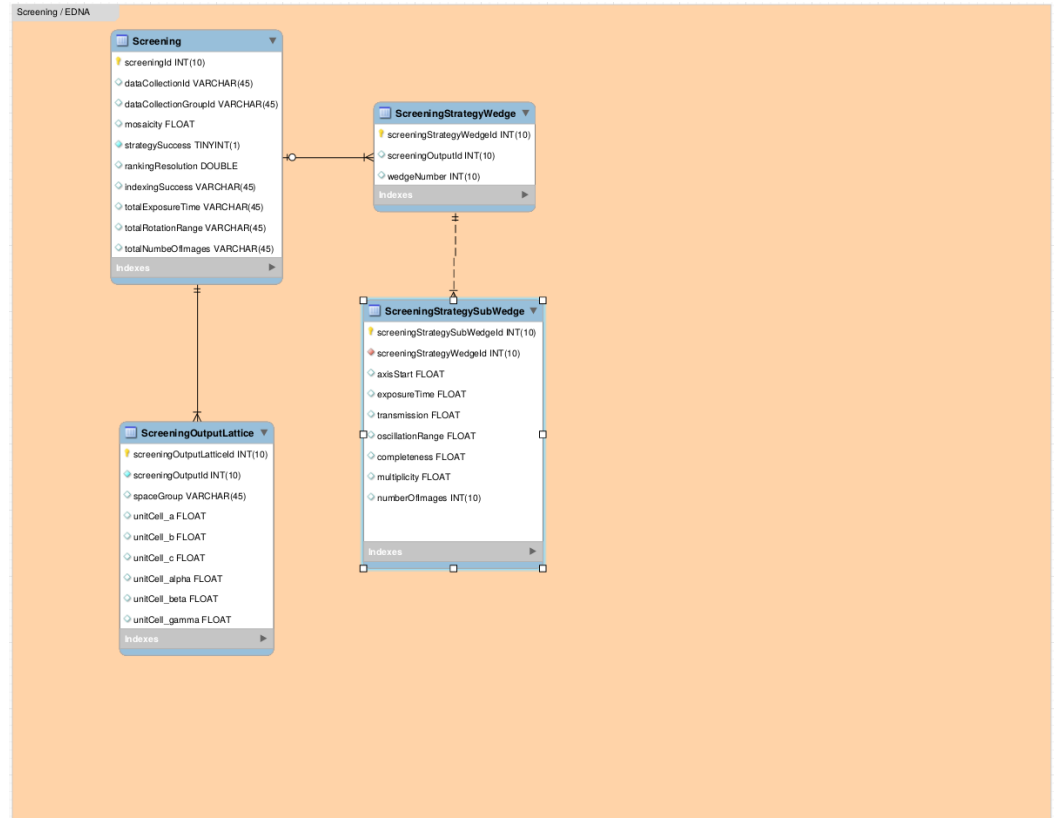
Screening Tables

- Proposal from @antonlinos
- Tables cover characterisation and strategy determination
- Intention to simplify and remove redundant columns
- Analysis of usage at ESRF/DLS...



Screening Tables

- ...Leads to vastly simpler layout
- Merge Screening, ScreeningOutput and Strategy tables, removing some fields
- Add an optional autoprocProgramID column
- Close to agreement on final version
- *Can we add twoTheta as variable to ScreeningStrategyWedge?*



Data Collection Tables

- Proposed by @karllevik
- Data collection, EnergyScan and XFEFlourescenceSpectrum share many columns
 - They are all in effect data collections so therefore should be represented as such in ISPyB
- A few specific columns from EnergyScan would be better stored as processing results
 - For example scanFileFullPath, jpegChoochFileFullPath etc. Can use more generic DataCollectionFileAttachments table

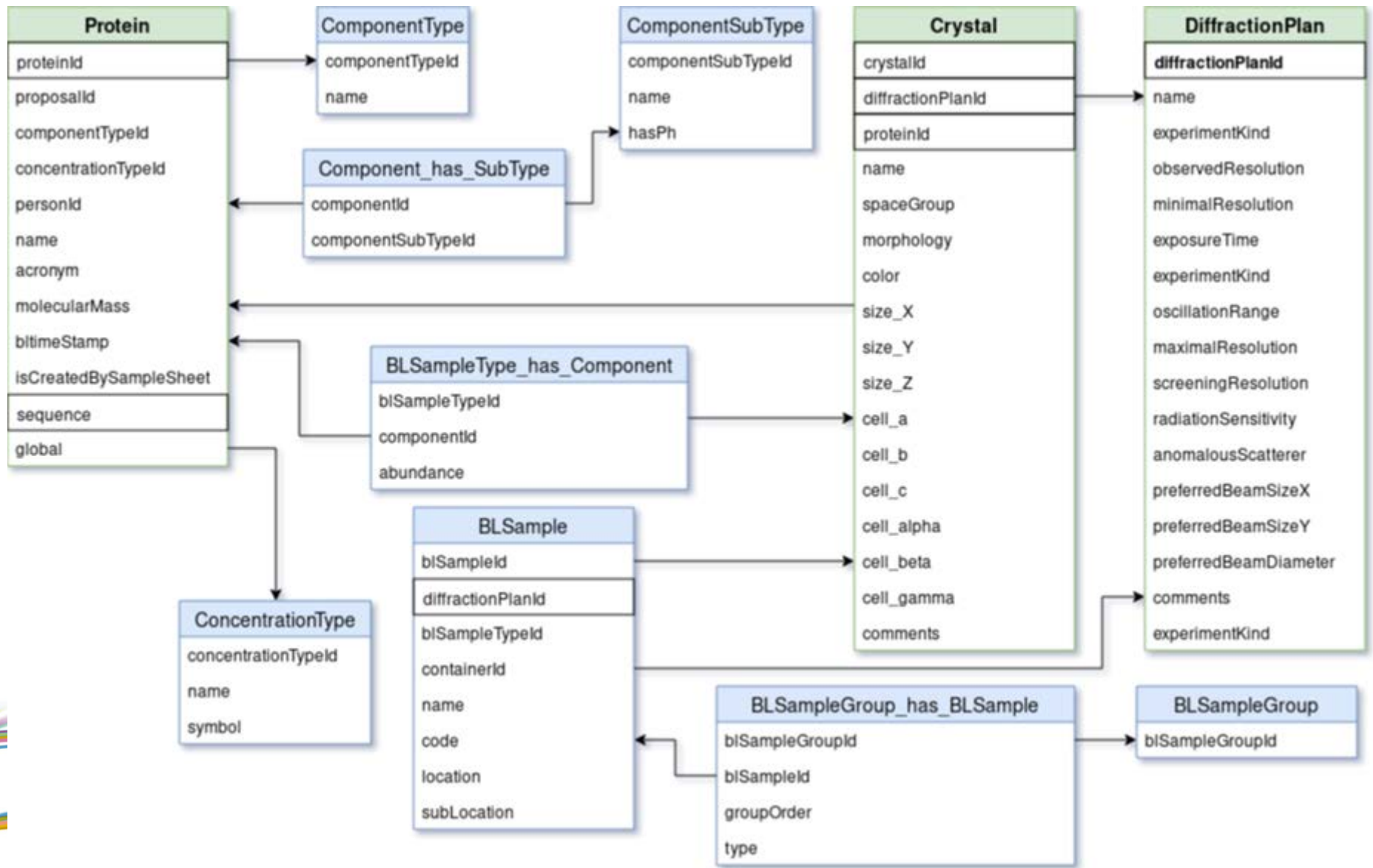
DC	EnergyScan	Fluorescence
dataCollectionId (PK)	energyScanId (PK)	xfeFluorescenceSpectrumId (PK)
DCG.sessionId	sessionId	sessionId
blSampleId	blSampleId	blSampleId
blSubSampleId	blSubSampleId	blSubSampleId
detectorId	fluorescenceDetector	
beamSizeAtSampleX	beamSizeHorizontal	beamSizeHorizontal
beamSizeAtSampleY	beamSizeVertical	beamSizeVertical
transmission	transmissionFactor	beamTransmission
comments	comments	comments
crystalClass	crystalClass	crystalClass
startTime	startTime	startTime
endTime	endTime	endTime
exposureTime	exposureTime	exposureTime
fileTemplate	filename	filename
averageTemperature	temperature	
wavelength		wavelength
totalAbsorbedDose or totalExposedDose?	xrayDose	
flux	flux	flux
flux_end	flux_end	flux_end
imageDirectory	workingDirectory	workingDirectory
axisStart,axisEnd		axisPosition

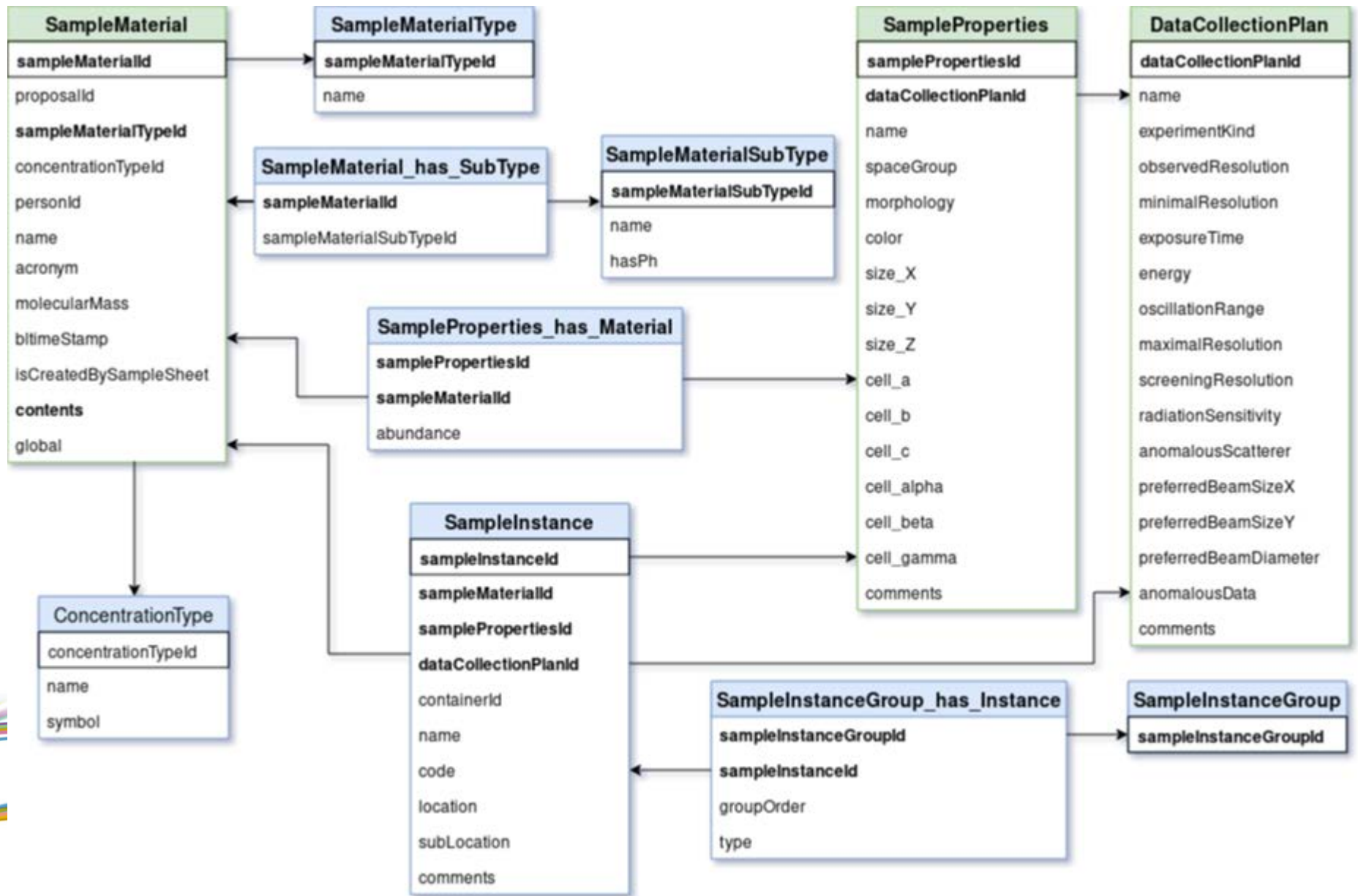
Data Collection Tables

- Agreement on initial proposal to merge tables into the Data Collection Table
- Wider discussion on whether DC table should be reduced to core set of values
- Currently DC is a "wide" table with many nulls
 - For example includes both MX and EM specific values
- Aesthetically an improved approach would see DC being a basic table with domain specific tables referring back to common Table
 - MXDataCollection, EMDataCollection
- From a performance perspective however, a wide table with NULLs is not an issue

Sample Tables

- Proposed by @karllevik
- Terminology from MX & life science
- Even with the realm of MX/EM/BioSAXS sample types vary
 - RNA, DNA, Protein, Virus
- Sample tables assume Protein/Crystal:
 - MX bias
 - BioSAXS tables are not well integrated with the rest of the schema (although scientifically relevant to their domain)





In Summary

- Renaming
 - Protein → SampleMaterial
 - Crystal → SampleProperties
 - BLSample → SampleInstance
 - BLSampleGroup → SampleInstanceGroup
 - ComponentType → SampleMaterialType
 - ComponentSubType → SampleMaterialSubType
 - DiffractionPlan → DataCollectionPlan
- ...and all the associated keys (proteinId => sampleMaterialId etc)
- However, without updating column names this leads to work without the payoff
- SampleProperties would still be a Crystal unless we add many fields
 - wide table problem again...
- Also does not accommodate BioSAXS use case

Sample Table Results

- More work required to develop a change that adds value and at least supports MX/EM and BioSAXS consistently
 - That should then provide better foundation to support other techniques
- Solution needs to address the distinction between sample components that are
 - the main component (e.g. protein) under investigation
 - Substances used as additives etc.
- Some consideration of reusing NXSample but initial thoughts are it does not fit well

Way forward

- Successful workshop
 - Some key agreements on refactoring Screening and Data Collection
- Appreciate resources limited
 - Refactoring to support future ISPyB requirements needs to be considerate of the impact
- Short term (12 months):
 - For samples - propose new tables that link to existing tables rather than other direction (non-invasive)
- Medium Term (12-24 months)
 - Seek closer alignment and consistency between MX/EM/BioSAXS
- Longer term (24 month+)
 - Should we have a roadmap/vision for what ISPyB is 2022?
 - Move towards micro services to aid reuse between sites?

Acknowledgements

Daniel Sanchez, ALBA
James Hall, Karl Levik, Diamond
Maxime Chaillet, Solange Delageniere, ESRF
Alejandro de Maria, Olof Svensson, ESRF
Stu Fisher, DLS/ESRF
Rasmus Fogh, Global Phasing
Idrissou Chado, Tatiana Isabet, SOLEIL

Thank you!