# Data and biophysics

Osvaldo Burastero
ARISE Fellow, Garcia-Alai Team
MOSBRI Course - Quality control for Integral Membrane Proteins
2022
14 September 2022
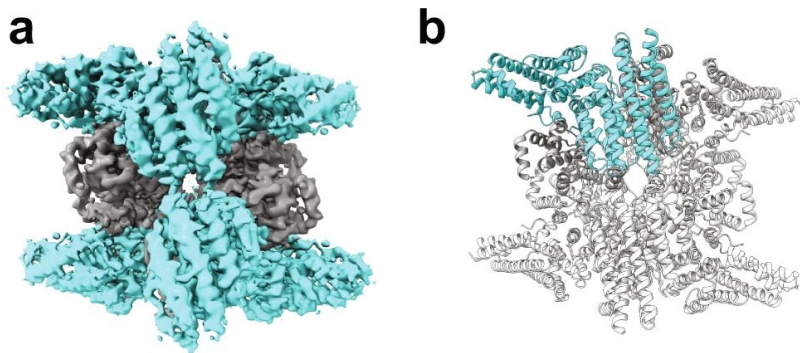
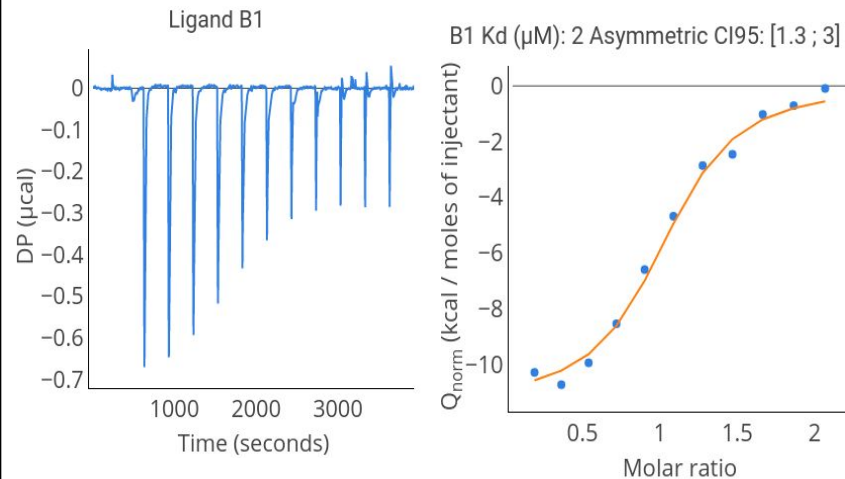# The final objective



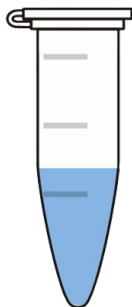| Cryo-EM structure of a 16-mer AENTH complex | Isothermal titration calorimetry example |
|---|---|

# How to get there



| 1. Experiment planning | 2. Sample preparation | 3. Quality check | 4. Measure & analysis |
|---|---|---|---|
| ● Data management plan <br> ● Predict possible outcomes | ● Protein expression & purification | ● Homogeneity <br> ● Integrity <br> ● Identity | ● Activity / inhibition assays <br> ● Structural studies |

# How to get there



| 1. Experiment planning | 2. Sample preparation | 3. Quality check | 4. Measure & analysis |
|---|---|---|---|
| <ul><li>Data management plan</li><li>Predict possible outcomes</li></ul> | <ul><li>Protein expression & purification</li></ul> | <ul><li>Homogeneity</li><li>Integrity</li><li>Identity</li></ul> | <ul><li>Activity / inhibition assays</li><li>Structural studies</li></ul> |

# How to get there



| 1. Experiment planning | 2. Sample preparation | 3. Quality check | 4. Measure & analysis |
|---|---|---|---|
| • Data management plan<br>• Predict possible outcomes | • Protein expression & purification | • Homogeneity<br>• Integrity<br>• Identity | • Activity / inhibition assays<br>• Structural studies |

# What is research data?

| Data | Metadata |
|---|---|
| <ul><li>Digital data generated during and after the research project<br><ul><li>Observations</li><li>Acquired data (raw & processed): text files, videos, …</li><li>Software (code, algorithms)</li></ul></li></ul> | <ul><li>Data about the data (provides context)<br><ul><li>Origin of the data</li><li>Who, when, why, how</li><li>Used resources</li><li>Licenses</li></ul></li></ul> |

MOSBRI
Molecular-Scale Biophysics
Research Infrastructure
EMBL

"Good metadata will save us precious time in the future"

# Data management plan (DMP)

| What is a DMP? | Why we need it? |
|---|---|
| ● Formal document that describes how the data will be handled during and after the project | ● Good scientific practice<br>● Required by funders and/or institution |



Horizon 2020
European Union Funding
for Research & Innovation

# Data management plan (DMP)

| What is a DMP? | Why we need it? | Benefits | FAIR Principles |
|---|---|---|---|
| ● Formal document that describes how the data will be handled during and after the project | ● Good scientific practice<br>● Required by funders and/or institution | ● Save time and resources<br>● Improved reproducibility and reusability | ● Findable<br>● Accesible<br>● Interoperable<br>● Reusable |

Horizon 2020
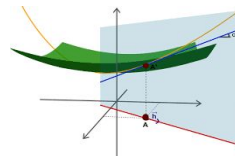European Union Funding
for Research & Innovation

scientific **data**

Explore content ⌄   About the journal ⌄   Publish with us ⌄

Open Access  |  Published: 15 March 2016

**The FAIR Guiding Principles for scientific data management and stewardship**

MOSBRI
Molecular-Scale Biophysics
Research Infrastructure

EMBL

# DMP checklist

| | | |
|---|---|---|
| Project | What is the project? | |
| Data | What type, format and size does the used and produced data have? | |
| Interpretation | Which information is required to understand the data? | |
| Procedures | Which procedures will be used to create, process and quality control the (meta)data? | |

# DMP checklist

| | | |
|---|---|---|
| Documentation | How the data processing steps will be recorded? |  |
| Access | Are there any security or access control requirements? |  |
| Project end | What happens to the data after the project finishes? |  |

# DMP checklist

| | |
|---|---|
| **Documentation** | How the data processing steps will be recorded? |
| **Access** | Are there any security or access control requirements? |
| **Project end** | What happens to the data after the project finishes? |
| **Intellectual property (IP)** | How will be the IP managed? |
| **Responsibilities** | Who is responsible for which part of the data management? |

# DMP checklist - zoom in

| Interpretation | |
|---|---|
| • Can the data be read only with specific software?<br>• Where is the data documentation to be found? Lab information management system?<br>• How is data going to be documented?<br>    • Metadata, identifiers (of biological entities) & ontologies | |

MOSBRI
Molecular-Scale Biophysics
Research Infrastructure
EMBL

# DMP checklist - zoom in

| Interpretation | Procedures |
|---|---|
| • Can the data be read only with specific software?<br>• Where is the data documentation to be found? Lab information management system?<br>• How is data going to be documented?<br>     • Metadata, identifiers (of biological entities) & ontologies | • How will data and files be named and organised?<br>• How will changes be tracked and propagated?<br>     • How will metadata and provenance be preserved?<br>     • How will derived data be updated? |

# DMP checklist - zoom in

| Processing |
| --- |
| <ul><li>Manual data processing steps</li><li>Configuration parameters</li><li>Analysis versions</li><li>Scientific workflow management system</li><li>Open source software</li><li>etc.</li></ul> |

MOSBRI
Molecular-Scale Biophysics
Research Infrastructure

EMBL

# DMP in real life - a living document

| Data generated for analysis of protein X stability and homogeneity | | | | |
|---|---|---|---|---|
| **Dataset** | **Origin** | **Size** | **Format** | **Availability** |
| DSF data | measurements performed in a nDSF Prometheus ® (Nanothemper) | <100 MB | .xlsx .csv | Open |
| DLS data | measurements performed in a DynaPro ® Plate Reader ( Wyatt Technology) | <10 MB | .csv | Open |

# How to get there



| 1. Experiment planning | 2. Sample preparation | 3. Quality check | 4. Measure & analysis |
|---|---|---|---|
| ● Data management plan<br>● Predict possible outcomes | ● Protein expression & purification | ● Homogeneity<br>● Integrity<br>● Identity | ● Activity / inhibition assays<br>● Structural studies |

# Knowledge discovery

**7** Novel knowledge

**6** Evaluation & interpretation

**5** Modeling

**4** Transformation

**3** Preprocessing

**2** Data collection & selection

**1** Understand the domain and objective

binding mechanism ΔH, ΔG, ΔS

affinity $K_D$

stoichiometry N

Country   Person   Expenditure

Document

Products

Income

# Sample Preparation & Characterization (SPC) Facility

- Optimisation
- Quality control
- Characterisation (thermodynamics & kinetics)

MALDI TOF

Differential Scanning Fluorimetry

Circular Dichroism

Mass Photometry

Isothermal Titration Calorimetry

MicroScale Thermophoresis
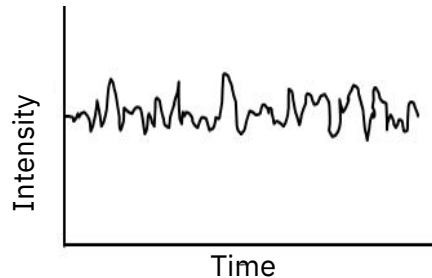
Dynamic Light Scattering

**and much more!**

# eSPC, enriching service provision

Differential Scanning
Fluorimetry (DSF)


FoldAffinity


MoltenProt

MicroScale
Thermophoresis
(MST)


ThermoAffinity

Mass Photometry
(MP)


PhotoMol

Burastero *et al.* (2022) / Acta Crystallogr. D

# Sample Preparation & Characterization (SPC) Facility

● Optimisation    ● Quality control    ● Characterisation (thermodynamics & kinetics)

MALDI TOF

Differential Scanning Fluorimetry

Circular Dichroism

Mass Photometry

Isothermal Titration Calorimetry

MicroScale Thermophoresis

Dynamic Light Scattering

**and much more!**

MOSBRI  Molecular-Scale Biophysics Research Infrastructure    EMBL

# Dynamic light scattering (DLS) in a nutshell

| Why? | How does it work? |
|---|---|
| <ul><li>Homogeneity of a sample</li><li>Estimation of the hydrodynamic radius (Hr)</li></ul> | <ul><li>It measures the autocorrelation of the scattered light</li></ul> |

# Dynamic light scattering (DLS) in a nutshell

| Why? | How does it work? | How do we extract valuable information? | Limitations |
|---|---|---|---|
| <ul><li>Homogeneity of a sample</li><li>Estimation of the hydrodynamic radius (Hr)</li></ul> | <ul><li>It measures the autocorrelation of the scattered light</li></ul> | <ul><li>We fit one/two, or a distribution of decay rates</li></ul> | <ul><li>Signal $\propto Hr^6$</li><li>Semi-quantitative</li><li>Scattering isn't isotropic for large particles</li></ul> |

# DLS - Theory

# DLS - Theory



**Small particles**
- Fast decay
- Give low signal
- Isotropic scattering

**Large particles**
- Slow decay
- Give high signal ($\propto Hr^6$)
- Anisotropic scattering

"Any model is at best, an useful fiction"

# DLS - Model & fitting

## What we measure

Second order correlation function $G_2$

$$G_2(\tau) = \langle I(t)I(t+\tau) \rangle$$

Integral over the product of intensities at time t and delayed time t+τ

$$g_2(\tau) = \frac{\langle I(t)I(t+\tau) \rangle}{\langle I(t) \rangle^2}$$

# DLS - Model & fitting

## What we measure

Second order correlation function $G_2$

$$G_2(\tau) = \langle I(t)I(t+\tau) \rangle$$

Integral over the product of intensities at time t and delayed time t+τ

$$g_2(\tau) = \frac{\langle I(t)I(t+\tau) \rangle}{\langle I(t) \rangle^2}$$

## Relationship with particle motion

Normalised first-order correlation function $g_1$

$$g_2(\tau) = 1 + \beta|g_1(\tau)|^2$$

Coherence factor $\beta \propto$ instrument & molecules

# DLS - Model & fitting

## What we measure

Second order correlation function $G_2$

$$G_2(\tau) = \langle I(t)I(t + \tau) \rangle$$

Integral over the product of intensities at time t and delayed time t+τ

$$g_2(\tau) = \frac{\langle I(t)I(t + \tau) \rangle}{\langle I(t) \rangle^2}$$

## Relationship with particle motion

Normalised first-order correlation function $g_1$

$$g_2(\tau) = 1 + \beta |g_1(\tau)|^2$$

Coherence factor β ∝ instrument & molecules

## Distribution of decay rates

Second order correlation function

$$g_1(\tau) = \int_0^\infty G(\Gamma) \exp(-\Gamma_\tau d\Gamma)$$

Intensity-weighted integral over a distribution of decay rates

# DLS - Model & fitting

## Decay rate & diffusion rates

Each decay rate can be associated to a certain diffusion factor $D$

$$D(s, q) = 1/(s * (q^2))$$

where s is the inverse of the decay rate and q is the Bragg wave vector ($\propto$ angle of detector & refractive index)

## Conversion to hydrodynamic radius

Sphere-like model allows estimating the Hr

$$Hr = \frac{k_b * temperature}{D * viscosity * 6\pi}$$

## The ISO recommended cumulants approach

- Moment analysis of the linear form of the measured correlogram

- Assumes a single particle family (Gaussian)

- Gives the Z-average (mean value) and the PdI (polydispersity index, relative variance of the Gaussian)

# DLS - Model & fitting

## Limitations of the cumulants approach

- Extremely sensitive to small amounts of aggregates

- Unsuitable for a polydisperse sample (polydispersity > 20 %)

## Fitting a distribution of decay rates

- We need to define a decay rate space

$$g_1(t) = \sum_{i=1}^{200} c_i \, exp \frac{-t}{s_i}$$

- Ill-posed problem that requires regularization

$$||Ax - b|| + \alpha||Mx|| + \beta||Ix||$$

$$\alpha \left|\left| \sum_{i=2}^{199} 2c_i - c_{i-1} - c_{i+1} \right|\right|$$

"DLS is (almost always) semi-quantitative"

# DLS - Fitting in practice



- Hr min value — 0.09 nm
- Hr max value — 5e5 nm
- Number of Hr points — 200
- Max time — 1 sec
- Alpha regularization — 0.1
- Beta regularization — 0

# DLS - Fitting in practice



- Hr min value — 0.09 nm
- Hr max value — 5e5 nm
- Number of Hr points — 200
- Max time — 1 sec
- Alpha regularization — 0
- Beta regularization — 0

# DLS - Fitting in practice



- Hr min value      0.09 nm
- Hr max value      5e5 nm
- Number of Hr points      200
- Max time      100 µs
- Alpha regularization      0.1
- Beta regularization      0

# DLS - Interpretation

| Good samples | Oligomerization | Hr values |
|---|---|---|
| • Cumulants PdI < 20 % (Malvern)<br>• One peak in region 1-20 nm with mass > 99.9 % | • Only big differences in size (factor 3-5) are detected, i.e., monomer to hexamer | • The Hr values are semi-quantitative<br>• Hr estimation assumes sphere-like model |

# Raynals, an app for DLS analysis



Step 1. Google "embl espc"



Step 2. Access spc.embl-hamburg.de



Step 3. Access Raynals spc.embl-hamburg.de/app/raynals

# Raynals, analysis code available under request

```python
W       = np.arange(1,len(data)+1) * 0 + 0.1 # all weights are equal, except the initial and last value
W       = W / np.max(W)
W       = np.append(W,np.array([1e2,1e2,1e2])) # weight to force the initial and last values equal to 0, and the sum of contributions equal to 1

rowToForceInitialValue      = np.zeros(kernel.shape[1])
rowToForceInitialValue[0] = 1
rowToForceLastValue         = np.flip(rowToForceInitialValue)

data    = np.sqrt(W) * np.append(data,np.array([1,0,0]))
data    = data.reshape(-1, 1)

kernel = np.vstack([kernel,np.ones(kernel.shape[1]),rowToForceInitialValue,rowToForceLastValue])
kernel = np.sqrt(W)[:, None] * kernel

cols    = kernel.shape[1]

M = np.zeros((cols,cols))
for i in range(1,M.shape[1]-1):
    M[i,i-1] = -1
    M[i,i]   =  2
    M[i,i+1] = -1

L       = alpha * M
C       = np.concatenate([kernel, L], axis=0)
d       = np.concatenate([data, np.zeros(cols).reshape(-1, 1)])

I       = beta * np.eye(*kernel.shape)
C       = np.concatenate([C, I], axis=0)
d       = np.concatenate([d, np.zeros_like(data)])
x, _    = nnls(C, d.flatten())

return x
```

# DLS - beyond the monodisperse / polydisperse sample

**Global Analysis of Dynamic Light Scattering Autocorrelation Functions**

Stephen W. Provencher*, Petr Štěpánek**

$$\hat{g}_1(t) = \int A(\tau)e^{-t/\tau}d\tau.$$

A($\tau$) is related to decay rates that can be converted into diffusion coefficients

$$\beta(q)\hat{g}_1(t;q) = \int_0^\infty A_r(\tau)e^{-t/\tau}d\tau + \int_0^\infty A_d(D)e^{-q^2 Dt}dD, \quad (2)$$

$$A(\tau;q) = A_r(\tau) + A_d[1/(q^2\tau)] \quad (3)$$



Decomposition into "diffusive" and "relaxational" (independent of angle) components

Provencher *et al.* (1996) / Part. Part. Syst. Charact.

# Sample Preparation & Characterization (SPC) Facility

- Optimisation
- Quality control
- Characterisation (thermodynamics & kinetics)

MALDI TOF

Differential Scanning Fluorimetry

Circular Dichroism

Mass Photometry

Isothermal Titration Calorimetry

MicroScale Thermophoresis

Dynamic Light Scattering

**and much more!**

# Mass photometry (MP) in a nutshell

| Why? | How does it work? |
|---|---|
| • Homogeneity of a sample <br> • Estimation of the molecular masses of different species | • Interference between scattered and reflected light combined with ratiometric imaging |

Soltermann *et al.* (2020) / Angew. Chem. Int. Ed

# Mass photometry (MP) in a nutshell

| Why? | How does it work? | How do we extract valuable information? | Limitations |
|---|---|---|---|
| • Homogeneity of a sample<br>• Estimation of the molecular masses of different species | • Interference between scattered and reflected light combined with ratiometric imaging | • We compare contrasts with known samples<br>• We fit $n$ distributions of masses | • nM concentration is required<br>• Detergent produces high background<br>• Accurate only with soluble proteins |



Incident light    Scattered light    Reflected light

Contrast    -0.01    0.01    1 μm

55 kDa    314 kDa    414 kDa    Counts    Mass (kDa)

refeyn.com
Soltermann *et al.* (2020) / Angew. Chem. Int. Ed

# MP - Theory

- Separation of incident from scattered and reflected light



Young *et al.* (2018) / Science
Cole *et al.* (2017) / ACS Photonics

# MP - Theory

| Experimental setup | Raw images |
|---|---|
| • Separation of incident from scattered and reflected light | • Interference: scattered - reflected<br>• High background |



Aqueous sample — Coverslip
Objective
Quarter-wave plate
Mirror — Polarizing beam splitter
Lens 1 — Telecentric lens 2
Lens 2 — Telecentric lens 1
Partial reflector
Lens 3 — Acousto-optic deflectors
445 nm laser
Camera

A
1300 ▬▬ 4000
Camera counts

Young *et al.* (2018) / Science
Cole *et al.* (2017) / ACS Photonics

# MP - Theory

| Experimental setup | Raw images | Preprocessing |
|---|---|---|
| • Separation of incident from scattered and reflected light | • Interference: scattered - reflected<br>• High background | • Background is removed by conversion to ratiometric images |

Young *et al.* (2018) / Science
Cole *et al.* (2017) / ACS Photonics

# MP - Theory

- Appear as a (dark) point spread function (PSF)
- The number should be >> unbinding events (bright spots)

# MP - Theory

| Binding events | Particle detection & quantification |
|---|---|
| • Appear as a (dark) point spread function (PSF) <br> • The number should be >> unbinding events (bright spots) | • Automated spot detection routine <br> • Fitting of candidate pixels (772 * 772 nm$^2$) to a 2D concentric gaussian model |



$$f(x,y) = A\left( e^{-\left[\frac{(x-x_0)^2}{2\sigma_x^2} + \frac{(y-y_0)^2}{2\sigma_y^2}\right]} - \frac{(1-T)}{s} e^{-\left[\frac{(x-x_0)^2}{2(s\sigma_x)^2} + \frac{(y-y_0)^2}{2(s\sigma_y)^2}\right]} \right) + b$$

$$A\left(1 - \frac{(1-T)}{s}\right)$$

Young *et al.* (2018) / Science

# MP - Theory

| Binding events | Particle detection & quantification | Contrast of a given particle |
|---|---|---|
| • Appear as a (dark) point spread function (PSF)<br><br>• The number should be >> unbinding events (bright spots) | • Automated spot detection routine<br><br>• Fitting of candidate pixels (772 * 772 nm$^2$) to a 2D concentric gaussian model | • Fitted from the contrast as a function of time |



$$f(x,y) = A\left(e^{-\left[\frac{(x-x_0)^2}{2\sigma_x^2} + \frac{(y-y_0)^2}{2\sigma_y^2}\right]} - \frac{(1-T)}{s} e^{-\left[\frac{(x-x_0)^2}{2(s\sigma_x)^2} + \frac{(y-y_0)^2}{2(s\sigma_y)^2}\right]}\right) + b$$

$$A\left(1 - \frac{(1-T)}{s}\right)$$

# MP - Model & fitting

Young *et al.* (2018) / Science

# MP - Model & fitting

## Contrast to molecular weight



## Calibration

- Using known protein standards
- Quantitative for soluble proteins without post translational modifications!



Young *et al.* (2018) / Science

# MP - Model & fitting

## Contrast to molecular weight



## Calibration

- Using known protein standards
- Quantitative for soluble proteins without post translational modifications!



## Unknown sample

- Fitting of gaussian distributions



Young *et al.* (2018) / Science

# PhotoMol, an app for MP analysis



Step 1. Google "embl espc"



Step 2. Access spc.embl-hamburg.de


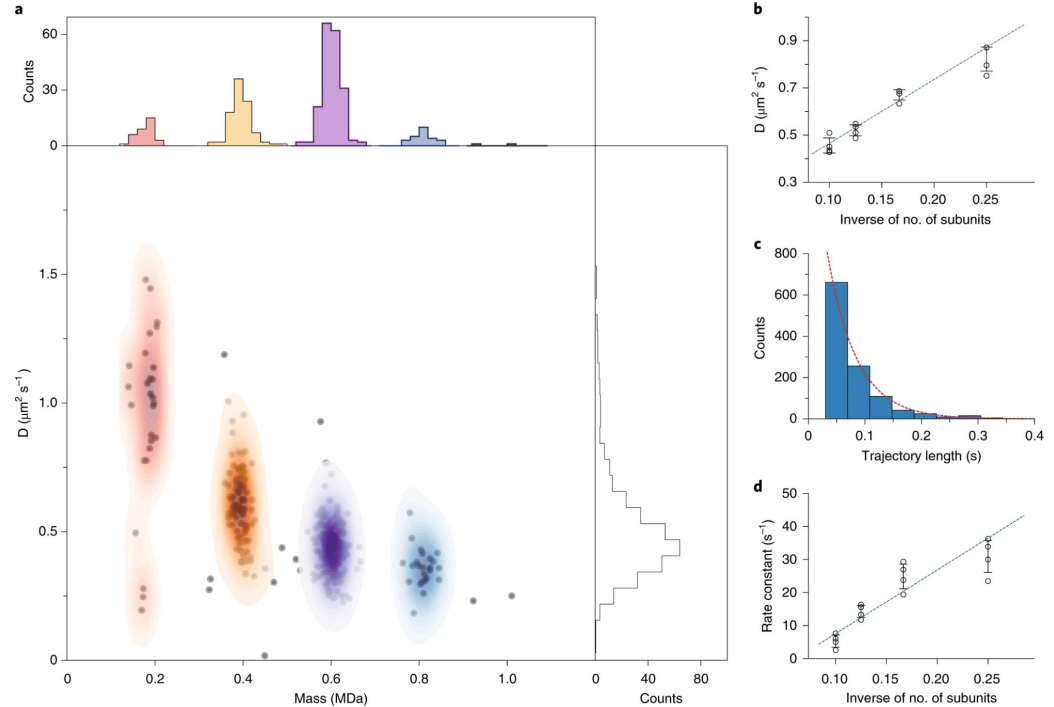
Step 3. Access PhotoMol spc.embl-hamburg.de/app/photoMol

# MP - beyond the yes/no complex formation question

**Mass photometry enables label-free tracking and mass measurement of single proteins on lipid bilayers**

Eric D. B. Foley, Manish S. Kushwah, Gavin Young & Philipp Kukura ✉

Simultaneously imaging, tracking and measuring the mass of diffusing biomolecular complexes on supported lipid bilayers (SLBs)

# MP - beyond the yes/no complex formation question

**Mass photometry enables label-free tracking and mass measurement of single proteins on lipid bilayers**

Eric D. B. Foley, Manish S. Kushwah, Gavin Young & Philipp Kukura ✉

Simultaneously imaging, tracking and measuring the mass of diffusing biomolecular complexes on supported lipid bilayers (SLBs)

# Sample Preparation & Characterization (SPC) Facility

- Optimisation
- Quality control
- Characterisation (thermodynamics & kinetics)

MALDI TOF

**Differential Scanning Fluorimetry**

Circular Dichroism

Mass Photometry

Isothermal Titration Calorimetry

MicroScale Thermophoresis

Dynamic Light Scattering

**and much more!**

# Nano Differential scanning fluorimetry (nDSF) in a nutshell

| Why? | How does it work? |
|---|---|
| • Protein stability or ligand binding | • By heating the sample and measuring the fluorescence |



Photo credit: Vadim Kotov

# Nano Differential scanning fluorimetry (nDSF) in a nutshell

| Why? | How does it work? | How do we extract valuable information? | Limitations |
|------|-------------------|----------------------------------------|-------------|
| • Protein stability or ligand binding | • Heating the sample and measuring the fluorescence | • Fitting unfolding models | • Autofluorescence<br>• Multiple transitions<br>• No transitions |



Native ⇌ Unfolded

# nDSF - Unfolding models

| n-mer | Protomers | States | Intermediates |
|-------|-----------|--------|---------------|
| Monomer | A | 2 | none |
| Monomer | A | 3 | 1 monomer |
| Monomer | A | 2+p | p monomers |
| Homodimer | A$_2$ | . | . |
| . | . | . | . |
| Heterodimer | AB | . | . |
| . | . | . | . |
| Heterotrimer | ABC | . | . |
| . | . | . | . |

Native → Unfolded

Native ⇋ Unfolded

Native ⇋ Intermediate ⇋ Unfolded

# nDSF - How to build a model

| | |
|---|---|
| Step 1. | Equation of the signal with *n* states |
| Step 2. | Pre/post transition dependence |
| Step 3. | States interconversion |

$Y_t(x) = Y_n[N] + \sum_j Y_j[I_j] + Y_u[U] + Y_d(x), \ j = 1, \ldots, p$

$Y_n = y_n + m_n x;$

Native ⇋ Unfolded

Send us code/equations & data and we will add them into the server!

# nDSF - The two state equilibrium model
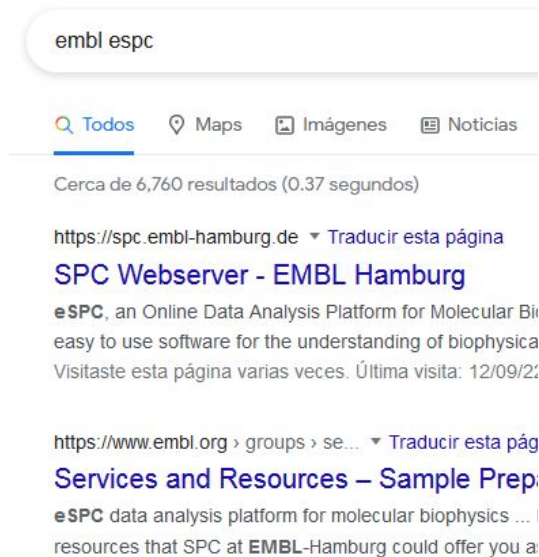
- Reversible equilibrium between native & unfolded states

$$\Delta G \left( T \right) = -RT \, \ln \left[ K \left( T \right) \right] = \Delta H_{\mathrm{m}} \left( 1 - T/T_{\mathrm{m}} \right) - \Delta C_{\mathrm{p}} \left[ \left( T_{\mathrm{m}} - T \right) + T \, \ln \left( T/T_{\mathrm{m}} \right) \right]$$

- ΔH and Tm can be precisely determined (but not ΔCp)

- ΔCp can be determined using different chemical denaturant concentrations

# MoltenProt & FoldAffinity, two apps for DSF analysis



Step 1. Google "embl espc"



Step 2. Access spc.embl-hamburg.de



Step 3. Access MoltenProt or FoldAffinity

SPC Team



EMBL
European Molecular Biology Laboratory

ARISE

MOSBRI
Molecular-Scale Biophysics
Research Infrastructure

EMBL Hamburg IT Team