

## OECD (Q)SAR Toolbox v.4.4.1

Example illustrating endpoint vs. endpoint  
correlation using ToxCast data

# Outlook

- **Background**
- Objectives
- The exercise
- Workflow

# Background

This presentation is designed to introduce the user to:

- ToxCast database as part of the Toolbox database
- Illustration of endpoint vs. endpoint correlations using:
  - ToxCast data
  - ToxCast and Estrogen receptor data

# Outlook

- Background
- **Objectives**
- The exercise
- Workflow

# Objectives

- This presentation demonstrates endpoint vs. endpoint correlations using ToxCast and Estrogen receptor data

# Outlook

- Background
- Objectives
- **The exercise**
- Workflow

## The exercise

- Illustration of endpoint data correlations using the ToxCast and estrogen binding data between the two types of data:
  - AC50 vs. AC50 endpoints associated with different test type
  - AC50 vs. Estrogen receptor binding data

# Outlook

- Background
- Objectives
- The exercise
- **Workflow**



# Workflow

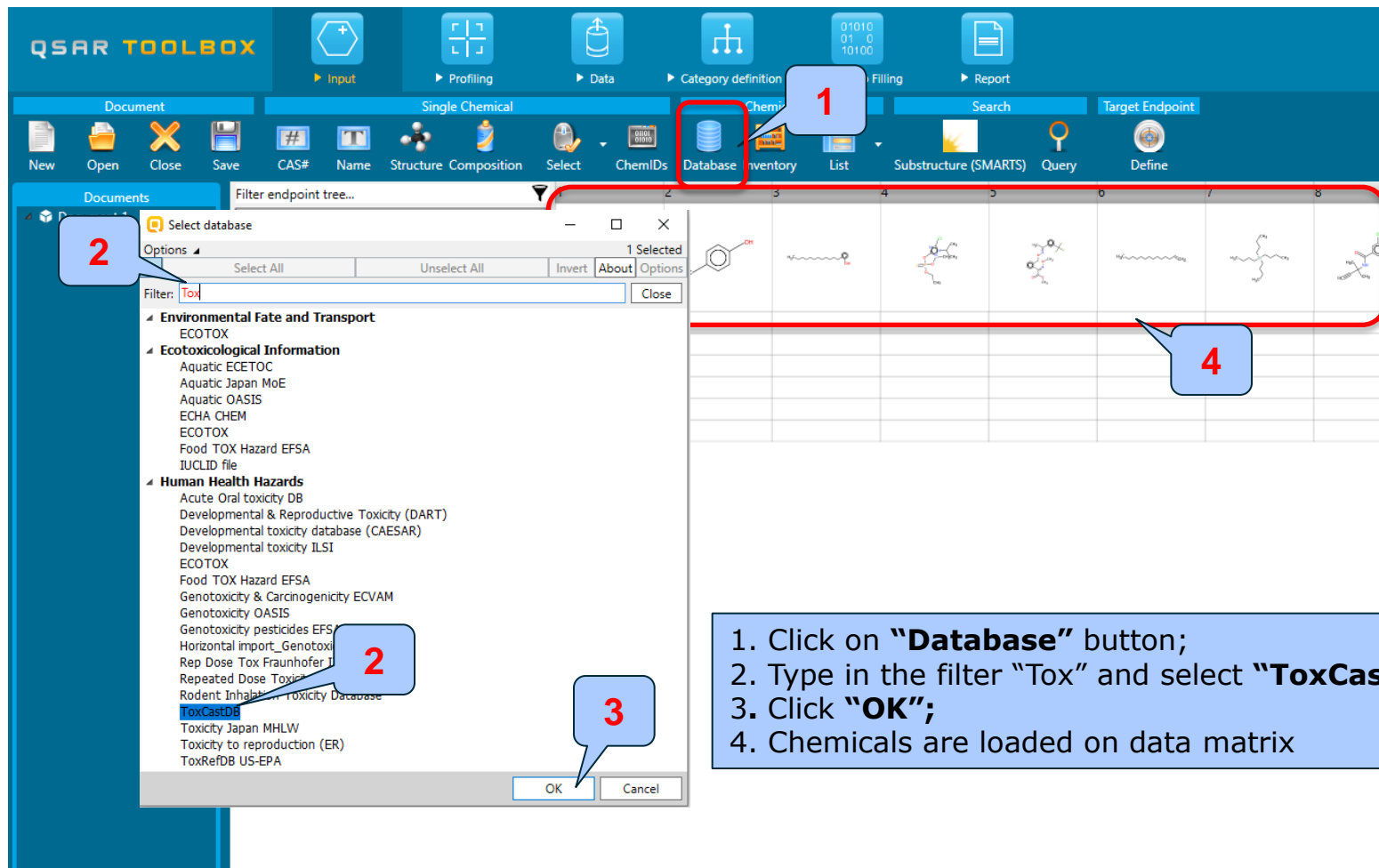
- **The Toolbox has six modules which are typically used in a workflow:**
  - Chemical Input
  - Profiling
  - Endpoints
  - Category Definition
  - Filling Data Gaps
  - Report
- **In this example we will use the modules in a different order, tailored to the aims of the example.**

# Outlook

- Background
- Objectives
- The exercise
- **Workflow**
  - **Load ToxCast database**

# ToxCast database

## Loading database



# ToxCast database

## Sidebar of database relevancy

Once the endpoint is selected, the relevant databases become highlighted in green.

For the purpose of forthcoming example we need to collect data for chemicals.

1. Expand **Human health hazard** and click on the level **ToxCast**;

2. The database is getting green highlighted; select it and click **Gather** (3);

4. Click **OK** to extract data from selected database.

# ToxCast database

## Data gathering

Filter endpoint tree...

Structure

Parameters

- Physical Chemical Properties
- Environmental Fate and Transport
- Ecotoxicological Information
- Human Health Hazards
  - Acute Toxicity
  - ADME
  - Bioaccumulation
  - Carcinogenicity
  - Developmental Toxicity / Teratog...
  - Genetic Toxicity
  - Immunotoxicity
  - Irritation / Corrosion
  - Neurotoxicity
  - Photoinduced toxicity
  - Repeated Dose Toxicity
  - Sensitisation
- ToxCast
  - ACEA 600/660
  - Apredica 425/2653
  - Attagene 1374/11710
  - BioSeek 971/21906
  - NCGC 1475/6890
  - Novascreen 975/8054
  - Odyssey Thera 969/2794
  - Undefined Assay provider 2/2
- Toxicity to Reproduction
- Toxicokinetics, Metabolism and D...

	1	2	3	4	5	6	7	8	9	10	11
Structure											
Parameters											
Physical Chemical Properties											
Environmental Fate and Transport											
Ecotoxicological Information											
Human Health Hazards											
Acute Toxicity											
ADME											
Bioaccumulation											
Carcinogenicity											
Developmental Toxicity / Teratog...											
Genetic Toxicity											
Immunotoxicity											
Irritation / Corrosion											
Neurotoxicity											
Photoinduced toxicity											
Repeated Dose Toxicity											
Sensitisation											
ToxCast											
ACEA 600/660			M: 0.0601 mg/L		M: 7.06 mg/L		M: 6.8 mg/L	M: 2.84 mg/L			M: 0.0219 mg/L
Apredica 425/2653					M: 1.69 mg/L				M: 32.1 mg/L		
Attagene 1374/11710		M: 0.88 mg/L	M: 0.113 mg/L	M: 0.627 mg/L	M: 4.82 mg/L	M: 16.2 mg/L	M: 0.033 mg/L	M: 3.79 mg/L	M: 3.4 mg/L		
BioSeek 971/21906	M: 0.127 mg/L	M: 0.16 mg/L		M: 0.464 mg/L	M: 0.539 mg/L	M: 0.243 mg/L		M: 0.663 mg/L	M: 0.464 mg/L	M: 0.187 mg/L	
NCGC 1475/6890	M: 0.367 mg/L		M: 0.156 mg/L	M: 1.61 mg/L	M: 0.357 mg/L		M: 1.86 mg/L		M: 0.000358 mg/L	M: 0.0144 mg/L	M: 6.23 mg/L
Novascreen 975/8054		M: 2.43 mg/L		M: 0.0957 mg/L	M: 0.209 mg/L	M: 0.0122 mg/L	M: 8.61 mg/L	M: 0.0597 mg/L			
Odyssey Thera 969/2794		M: 6.89 mg/L	M: 0.121 mg/L	M: 9.54 mg/L	M: 0.592 mg/L		M: 17.1 mg/L	M: 14.1 mg/L		M: 6.03 mg/L	
Undefined Assay provider 2/2											
Toxicity to Reproduction											
Toxicokinetics, Metabolism and D...											

1. The data appears in the datamatrix under level "ToxCast"

# Outlook

- Background
- Objectives
- The exercise
- **Workflow**
  - Load ToxCast database
  - **ToxCast database - overview**

# ToxCast database

## Background

- A major part of EPA's CompTox research is the ToxCast™ project. ToxCast is a multi-year project launched in 2007 that uses automated chemical screening technologies (called "high-throughput screening assays") to expose living cells or isolated proteins to chemicals. The cells or proteins are then screened for changes in biological activity that may suggest potential toxic effects. These innovative methods have the potential to limit the number of required laboratory animal-based toxicity tests while quickly and efficiently screening large numbers of chemicals.
- ToxCast has evaluated over 2,000 chemicals from a broad range of sources including: industrial and consumer products, food additives, and potentially "green" chemicals that could be safer alternatives to existing chemicals. Chemicals were evaluated in over 700 high-throughput assays that cover a range of high-level cell responses and approximately 300 signaling pathways.
- ToxCast results are contributed to the federal agency collaboration called Toxicity Testing in the 21st Century (Tox21). Tox21 pools chemical research, data and screening tools from multiple federal agencies including the National Toxicology Program. So far, Tox21 has compiled high-throughput screening data on nearly ten thousand chemicals.

# Outlook

- Background
- Objectives
- The exercise
- **Workflow**
  - Load ToxCast database
  - ToxCast database – overview
  - **Correlation of data - background**



# Correlation of endpoint data

## Background

- This functionality introduces the user to the opportunity to analyze correlations between selected gap filling endpoints (endpoints used for prediction) and other endpoint data.
- It is applicable for correlation analysis of data presented in ordinary, interval or ratio scale.
- If correlated data are measured in interval or ratio scale they are transformed in ordinary scale and the strength of the correlation is estimated by Spearman correlation coefficient.
- Basically, this functionality provides a correlation between a target endpoint (this is the initial endpoint selected by the user) displayed on ordinate axis (Y-axis) and other endpoint data displayed on the abscissa (X-axis).

# Correlation of endpoint data

## Spearman coefficient factor

- Spearman's rank correlation coefficient is a nonparametric rank statistic proposed by Charles Spearman as a measure of the strength of an association between two variables. It assesses how well the relationship between two variables can be described using a monotonic function.
- Spearman correlation coefficient could be used for exploring the covary between:
  - two ranked variables
  - one measurement variable and one ranked variable (in this case, the measurement variable need to be to converted to ranks)
- Spearman correlation varies from -1 to +1 and the interpretation of the coefficient factor is provided below:
  - 0.00 – 0.19 – very weak correlation
  - 0.20 – 0.39 – weak correlation
  - 0.40 – 0.59 – moderate correlation
  - 0.60 – 0.79 – strong correlation
  - 0.80 – 1.0 – very strong

# Outlook

- Background
- Objectives
- The exercise
- **Workflow**
  - Load ToxCast database
  - ToxCast database – overview
  - Correlation of data – background
  - **Types endpoint correlations**

## Types endpoint correlations

**Types endpoint correlations are as follows:**

- Continuous vs. continuous
- Categorical vs. categorical\*:
  - ✓ Categorical vs. categorical
  - ✓ Categorized continuous vs. categorical
  - ✓ Categorized continuous vs. categorized continuous

\*All type categorical vs. categorical correlations are not illustrated in this presentations. These type correlations are shown in presentation "Tutorial 13 TB 4.4.1 Example illustrating endpoint vs. endpoint correlation for apical endpoints"

# Outlook

- Background
- Objectives
- The exercise
- **Workflow**
  - Load ToxCast database
  - ToxCast database – overview
  - Correlation of data – background
  - **Types endpoint correlations**
    - Continuous vs. continuous

# Types endpoint correlations

## Continuous vs. continuous

- The aim of this type correlation is to illustrate how continuous type endpoint data or so called ratio data correlate with each other (e.g. LC50 vs. EC50 data)
- In this example we will illustrate how AC50 data associated with two different test assays extracted from ToxCast DB correlate with each other:
  - NCGC Reporter Gene Assay ERa Agonist, Estrogen receptor 1 (assay 1)
  - Tox21\_Era\_BLA\_Agonist\_ch2 (assay 2)
- Step by step workflow is presented on the next few slides. Summary of the workflow steps are provided below:
  - *Gather experimental data (step 1)*
  - *Selection of target endpoint (step 2)*
  - *Enter Gap filling (step 3)*
  - *Change default X-descriptor (logKow) with AC50 data (step 4)*

# Types endpoint correlations

Continuous vs. continuous

*Gather experimental data – step 1*

The screenshot shows the QSAR TOOLBOX interface. The top toolbar has buttons for 'Data', 'Input', 'Profiling', 'Data Gap Filling', and 'Report'. The 'Data' button is highlighted with a red circle and the number 1. On the left, the 'Databases' list includes 'ToxCastDB' which is selected with a red circle and the number 2. Below the 'Databases' list, the 'Gather' button is highlighted with a red circle and the number 3. The main area shows a 'Filter endpoint tree...' on the left and a data matrix on the right. The data matrix has columns for chemical structures and rows for various endpoints like 'Acute Toxicity', 'ADME', 'Bioaccumulation', etc. The 'ToxCast' section is expanded, showing data for various assays.

Follow the steps if you do not already load Toxcast data on data matrix.

1. Go to **"Data"**
2. Select **"ToxCast"** DB
3. Click **"Gather"**

# Types endpoint correlations

Continuous vs. continuous

*Selection of target endpoint – step 2*

**3**

**1**

**2**

**4**

**Possible data inconsistency**

Metadata

- Assay
- Assay provider
- Endpoint
  - ☒ AC50 (374 chemicals; 505 data)
- Entrez gene name
  - ☒ estrogen receptor 1 (374 chemicals; 505 data)
- Native scale/unit
  - ☒  $\mu\text{M}$  (374 chemicals; 505 data)
- Test organisms (species)

Select scale/unit to use

- ☐  $\text{g/m}^3$  [0 native data and 505 converted]
- ☒  $\log(1/\text{mol/L})$  [0 native data and 505 converted]
- ☐  $\mu\text{M}$  [505 native data and 0 converted]

Converted data

505 from scale/unit  $\mu\text{M}$

Chemicals 374/374; Data 505/505

OK Cancel

In this step we need to select the first endpoint, which will be used in the correlation. This will be the endpoint displayed on the Y-axis. The second endpoint will be selected latter on.

1. Go to **Data Gap Filling** module;
2. Highlight the empty cell next to the AC50 endpoint associated with assay: "NCGC Reporter Gene Assay ERa Agonist"
3. Click **"Trend analysis"**;
4. A window alerting you for data inconsistencies appears. Keep it as it is. Click **"OK"**.

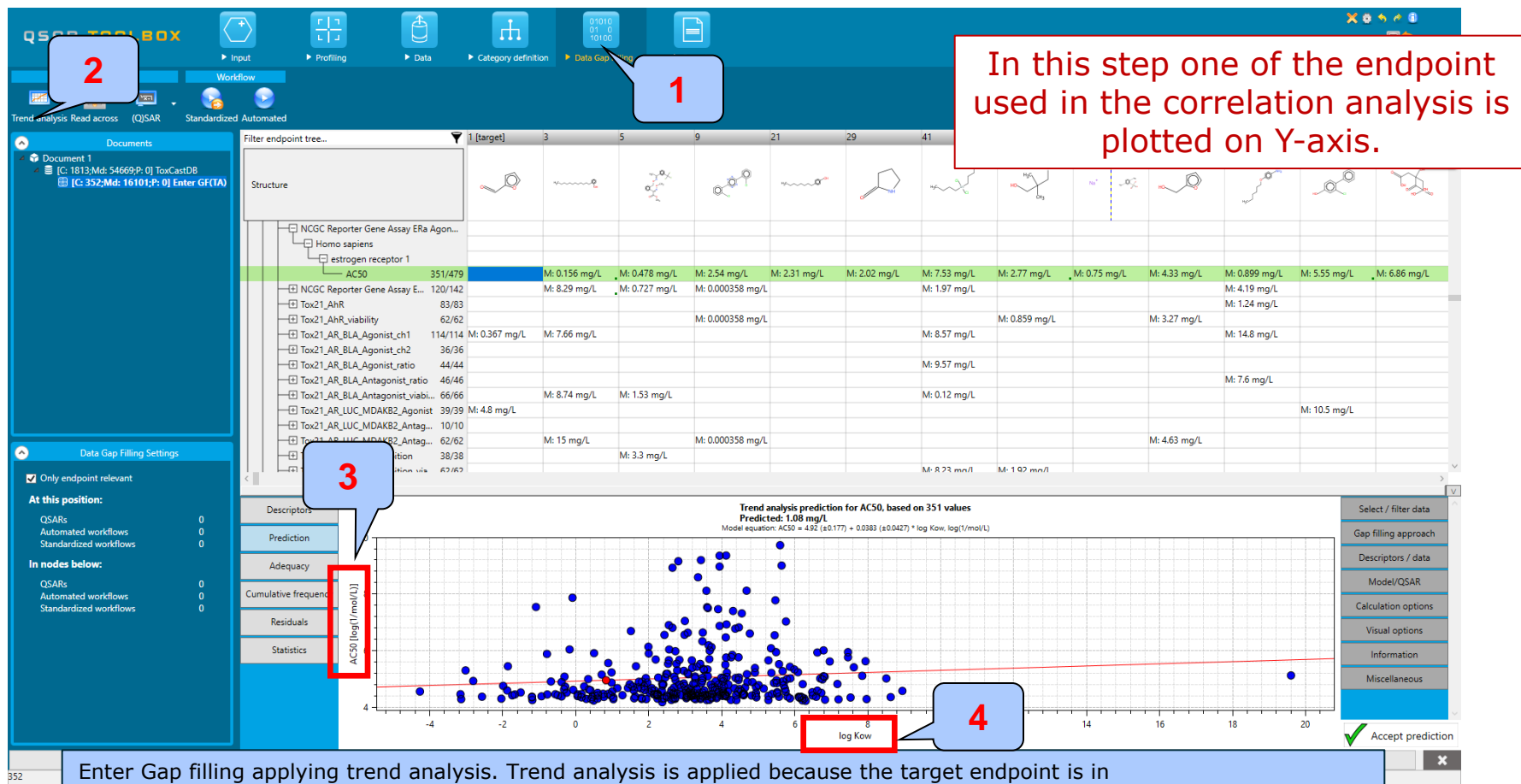




# Types endpoint correlations

Continuous vs. continuous

Enter Gap filling – step 3



# Types endpoint correlations

Continuous vs. continuous

Replacement of default X-descriptor (logKow) with AC50 data – step 4

In this step we need to select the second endpoint used for the correlation analysis. This endpoint will replace the default X- parameter (usually logKow)

3

1

2

Trend analysis prediction for AC50, based on 351 values  
Predicted: 1.08 mg/L  
Model equation:  $AC50 = 4.50 (\pm 0.177) + 0.0383 (\pm 0.0427) \cdot \log Kow, \log(1/mol/L)$

AC50 [log(1/mol/L)]

1. Click on "Descriptors / data";
2. Click "Select endpoint tree descriptor"; A window with arranged "Endpoint data tree" appears.
3. Expand Toxcast level to the level NCGC (352/2542).

# Types endpoint correlations

Continuous vs. continuous

Replacement of default X-descriptor (logKow) with AC50 data – step 4

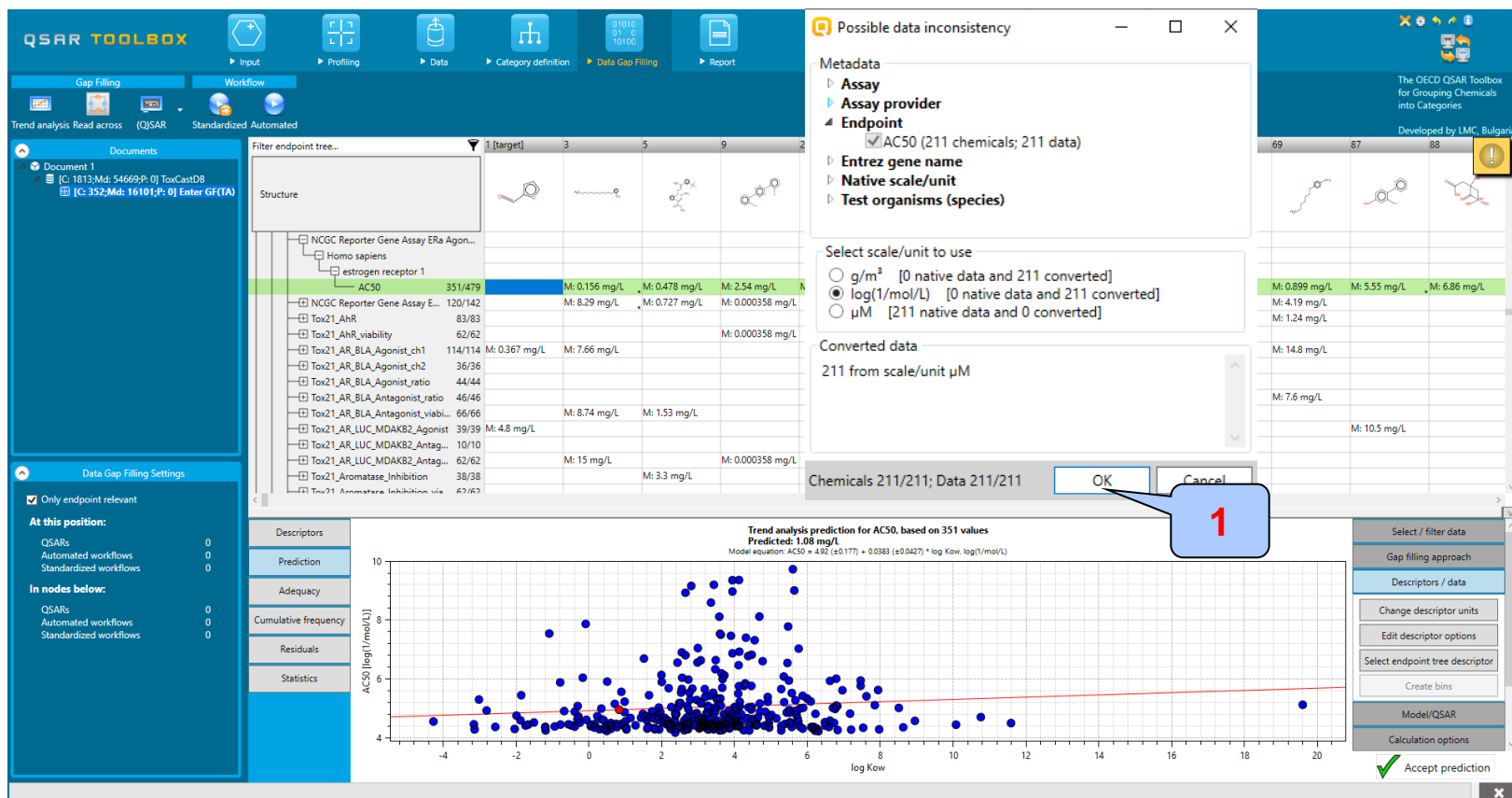
The screenshot displays the 'Select endpoint descriptor' dialog box in the QSAR Toolbox. The dialog lists various endpoints under the 'NCGC' node. A red box highlights the 'AC50 (211/211)' endpoint. A blue callout '1' points to the 'NCGC' node, a blue callout '2' points to the 'AC50 (211/211)' endpoint, and a blue callout '3' points to the 'OK' button. The background shows a chemical structure grid and a trend analysis plot for AC50.

1. Click on "NCGC" node to open the sub-nodes;
2. Select endpoint, which will be placed on X-axis circled in red box; point the mouse on the level of **AC50 (211/211)**;
3. Click "OK" button.

# Types endpoint correlations

Continuous vs. continuous

Replacement of default X-descriptor (logKow) with AC50 data – step 4

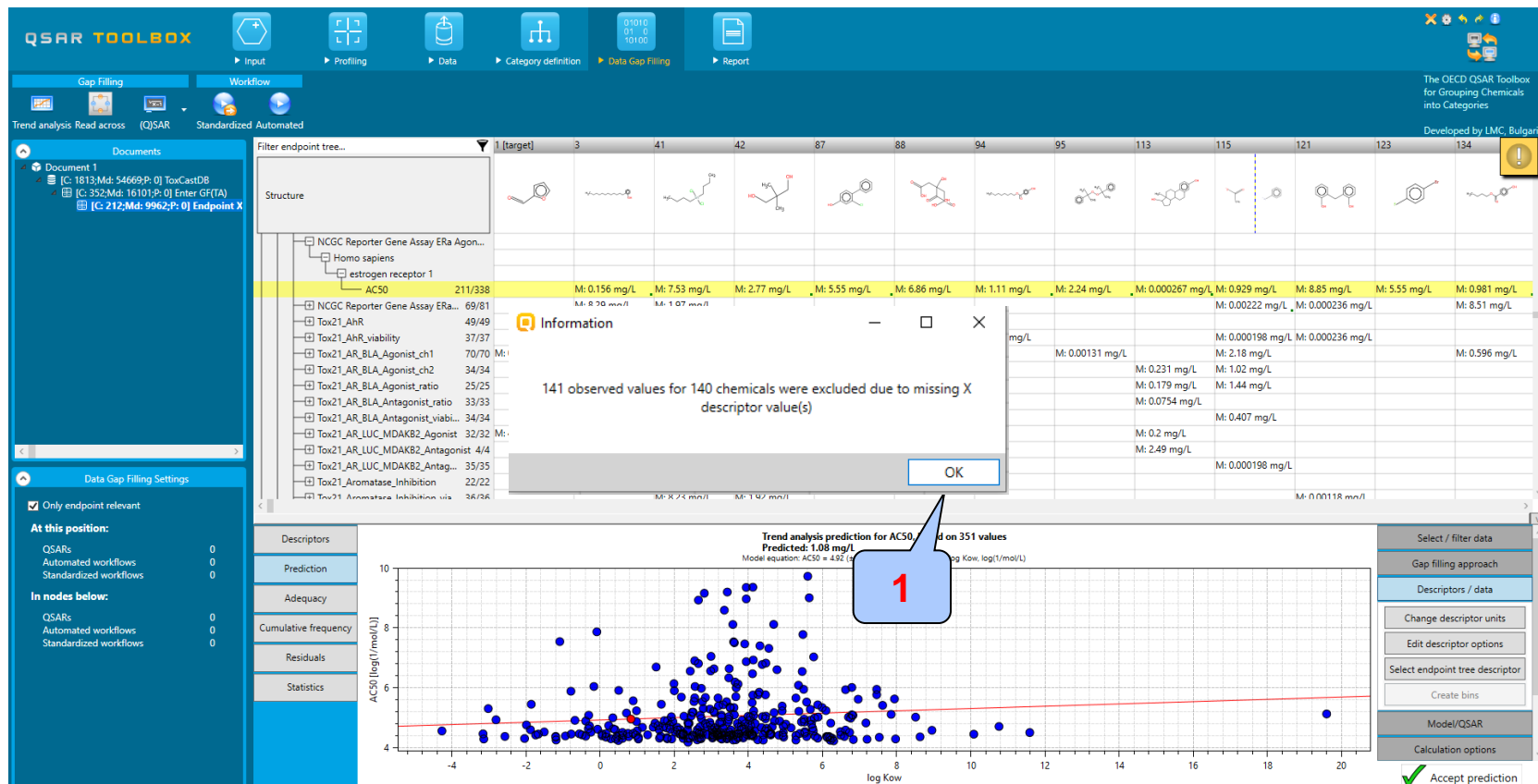


1. Click "OK" on the message alerting you for data inconsistency; The aim of this example is to see how the data correlates.

# Types endpoint correlations

Continuous vs. continuous

Replacement of default X-descriptor (logKow) with other AC50 data – step 4

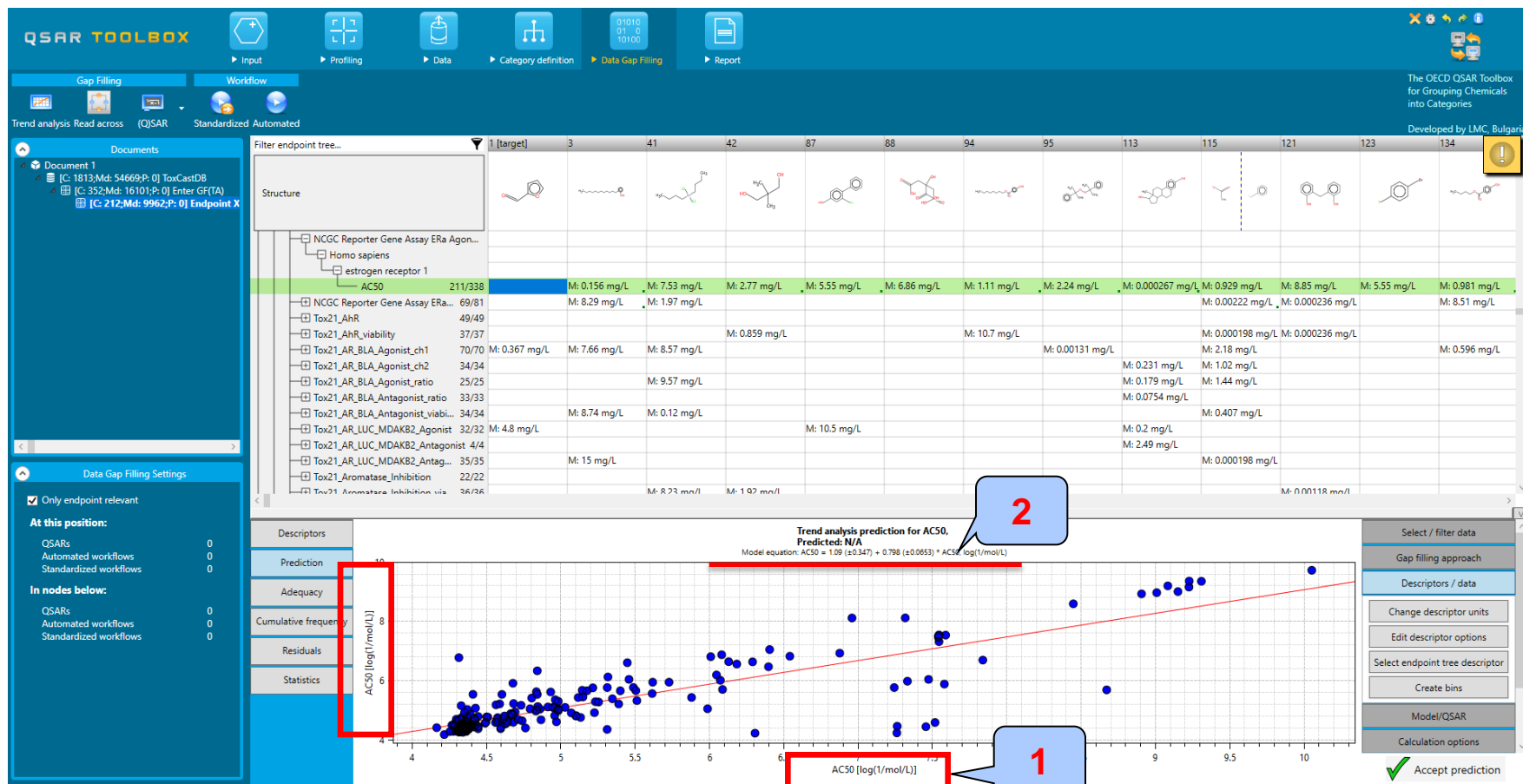


1. Click "OK" on the message informing you for the number of excluded chemicals due to missing X-descriptor data. They are analogues which do not have AC50 data for the assay "Tox21...", plotted on X-axis. This will not affect the value of correlation coefficient;

# Types endpoint correlations

Continuous vs. continuous

Replacement of default X-descriptor (logKow) with other AC50 data – step 4



1. The graph obtained after replacing of log Kow with Toxcast endpoint is visualized;
2. The equation including endpoint data is rebuild;

# Types endpoint correlations

Continuous vs. continuous

*Interpretation of correlation results*

- In this example, we have correlated two AC50 endpoints associated with different type assay
- As seen from the graph, a linear relationship between two endpoints has been observed
- In order to assess only the chemicals having positive estrogen binding activity we remove the “Non-binders” chemicals based on subcategorization by “Estrogen receptor binding by OASIS” profiler (illustrated on next slide)



# Types endpoint correlations

## Continuous vs. continuous

### Subcategorization by Estrogen receptor binding profiler

In this stage chemicals which do not shows estrogen binding activity will be eliminated.

1. Open "Select/filter data" menu item, then click "Subcategorize";

2. Select "Estrogen receptor binding" profiler;

3. Select only "Non binder" categories by left mouse click and hold "Ctrl" button;

4. Click "Remove" button.

# Types endpoint correlations

Continuous vs. continuous

*Interpretation of correlation results*

In the forthcoming slides are illustrated three endpoint vs. endpoint correlations:

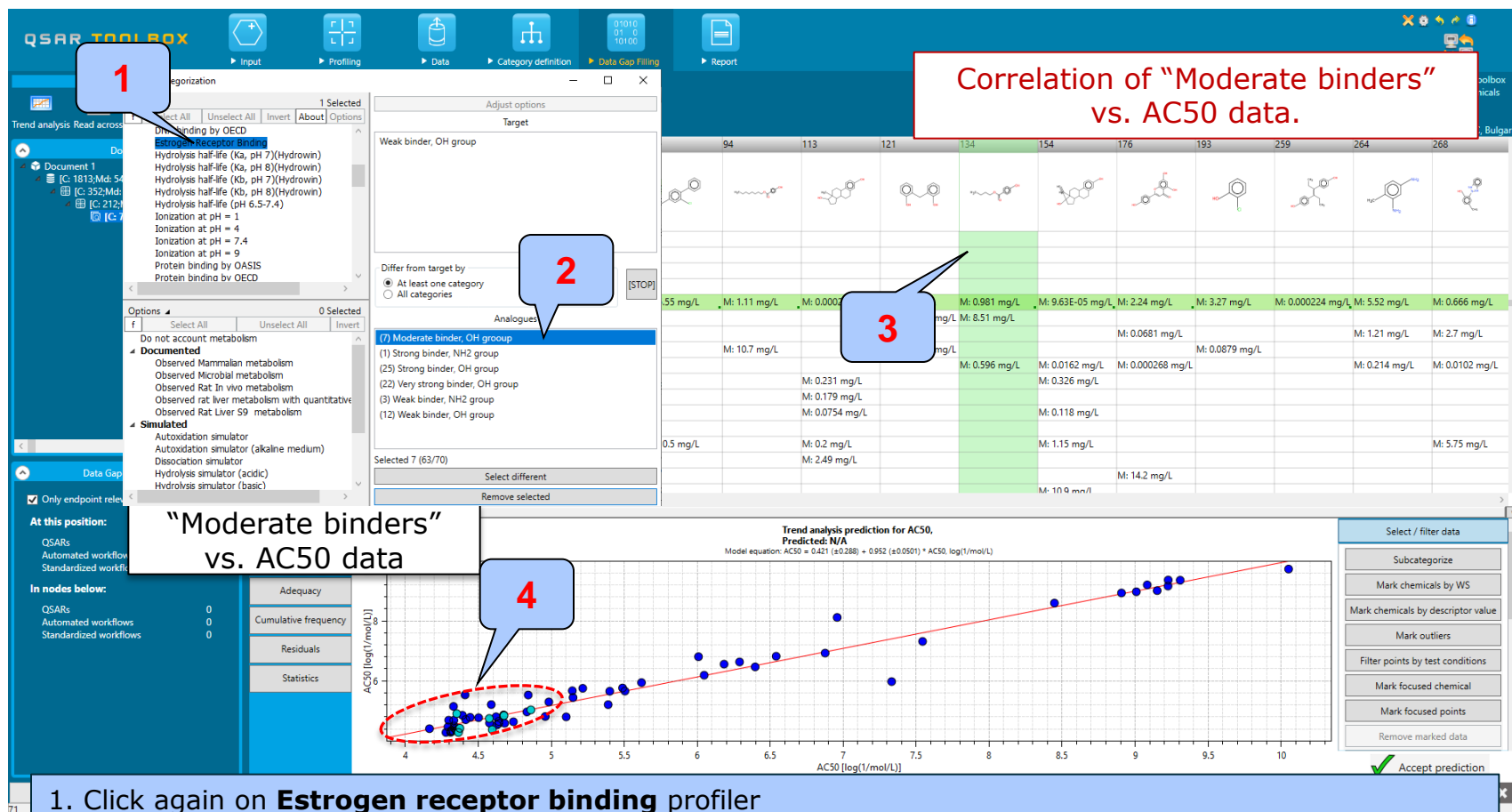
- Correlation of "Moderate ER binders" vs. AC50 data;
- Correlation of "Weak ER binders" vs. AC50 data;
- Correlation of "Strong ER binders" vs. AC50 data.

The aim of the slides is to illustrate how the chemicals possessing ER binding potency correlate with AC50 data.

# Types endpoint correlations

Continuous vs. continuous

*Subcategorization by Estrogen receptor binding profiler*

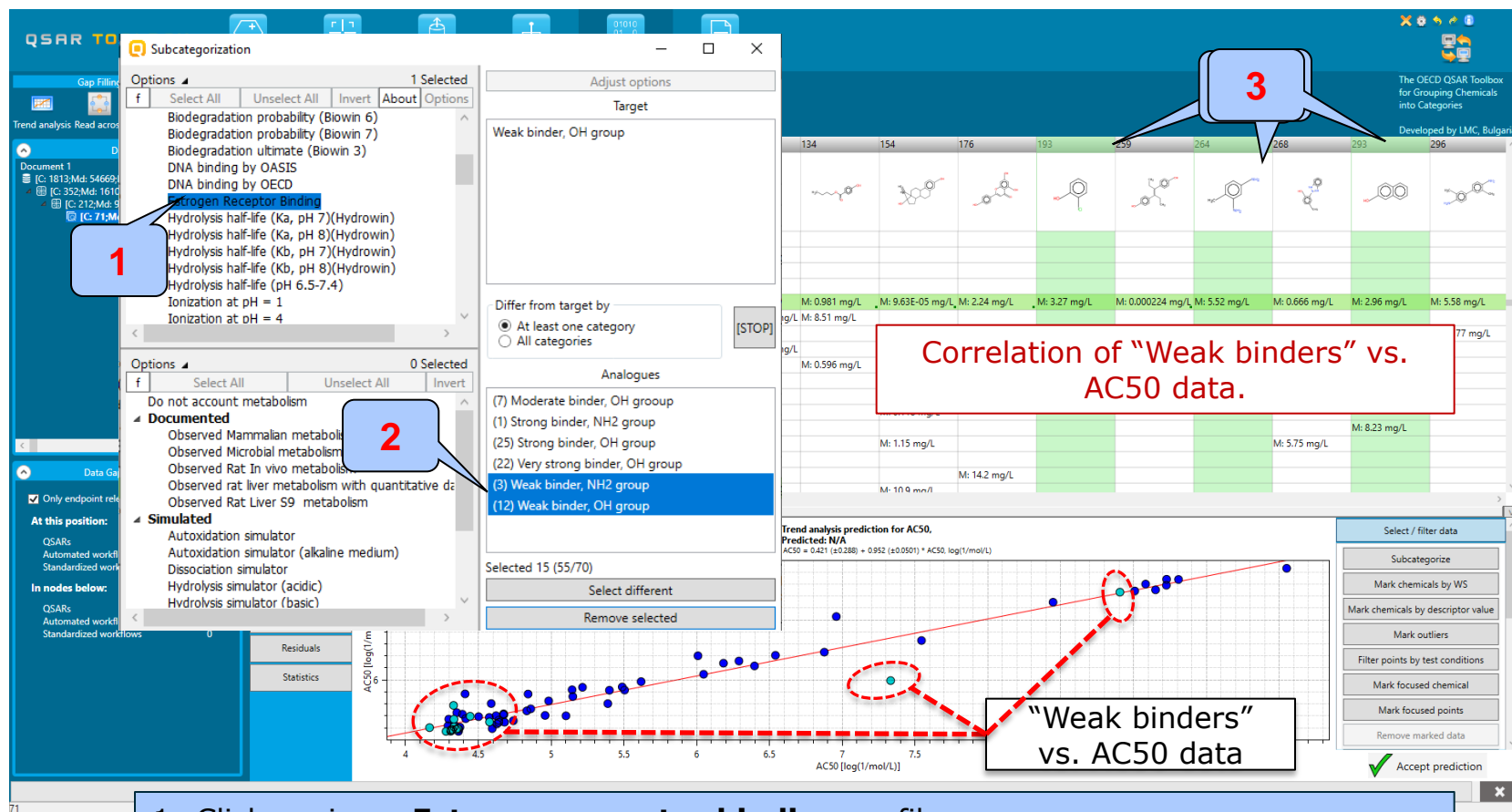


1. Click again on **Estrogen receptor binding** profiler
2. Select **"Moderate binder"** categories
3. The chemicals corresponding to the selected categories are highlighted in green on data matrix and in light blue on the graph (4)

# Types endpoint correlations

## Continuous vs. continuous

### Subcategorization by Estrogen receptor binding profiler

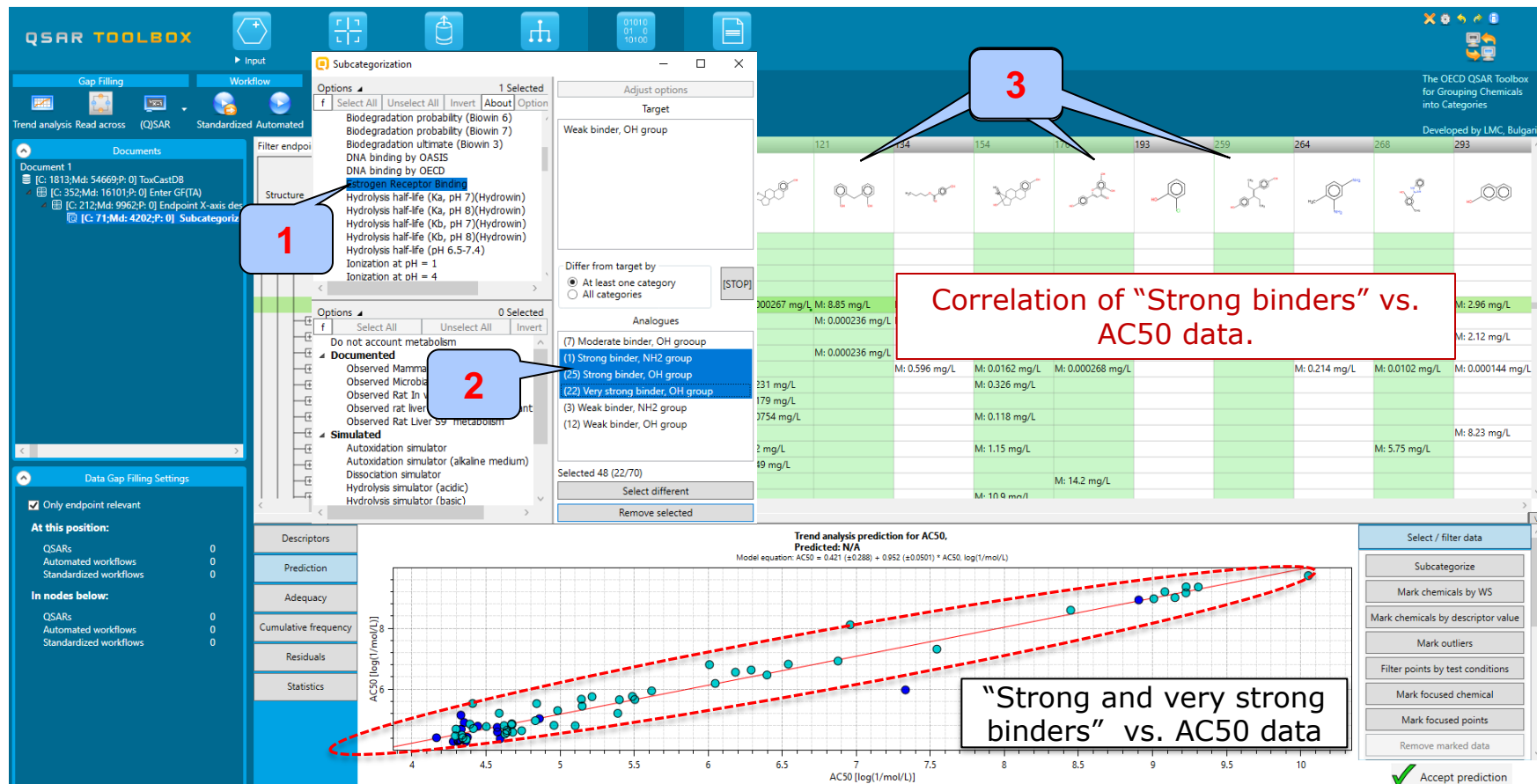


1. Click again on **Estrogen receptor binding** profiler
2. Select **"Weak binder"** categories (left mouse click and hold "Ctrl" button);
3. The chemicals corresponding to the selected categories are highlighted in green;

# Types endpoint correlations

## Continuous vs. continuous

### Subcategorization by Estrogen receptor binding profiler



# Types endpoint correlations

Continuous vs. continuous

*Correlation results*

- The two AC50 endpoints associated with different types of assays have been correlated each other
- Non binders according to the Estrogen receptor binding profiler have been eliminated from the correlation
- User can analyse the distribution of remaining ER binders (Very strong, Strong, Moderate and Weak) across selected AC50 endpoint