# A Proposition of XML Format for Proteomics Database

Ken'ichi KAMIJO, Toshimasa YAMAZAKI, and Akira TSUGITA

Proteomics Research Center,

Fundamental Research Labs., NEC Corp.

# Data Format Standardization

- **Download entries from public DBs as a flat-file**
  - easy for a person to read
  - different formats for every DB
  - sometimes needs special access methods and special applications for each format

- **Needs machine-readable formats for software tools**

- **To boost studies by exchanging data among researchers**

Activates standardization

# XML format

- **XML (eXtensible Markup Language)**
  - Highly readable for machine and person
  - Can represent information hierarchy and relationships
  - Details can be added right away
- **Convenient for exchanging data**
  - Easy to translate to other formats
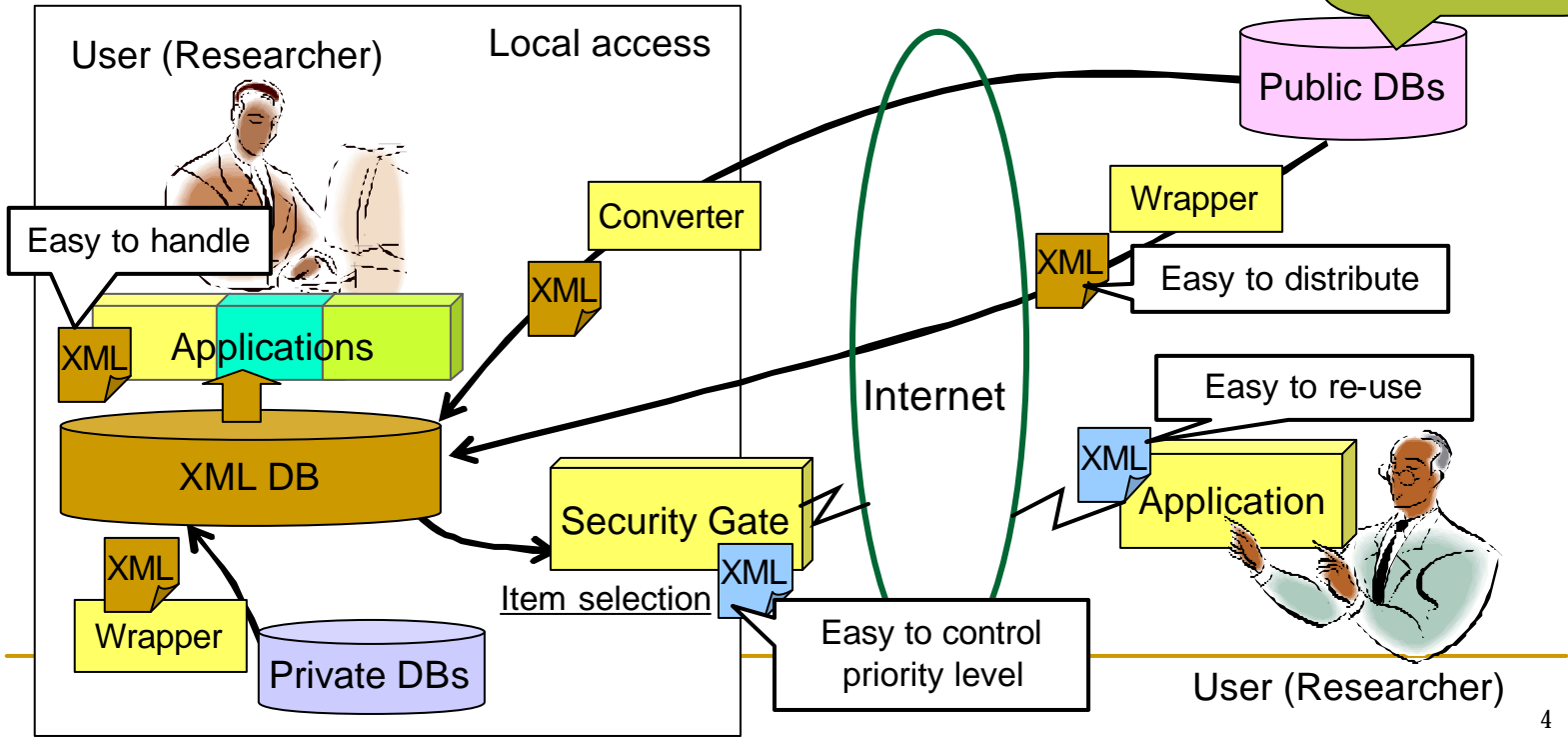  - Logical-check by a Document Type Definition (DTD)

```
<tag_source  element_growth="8 weeks">
 rice leaf
</tag_source>
```
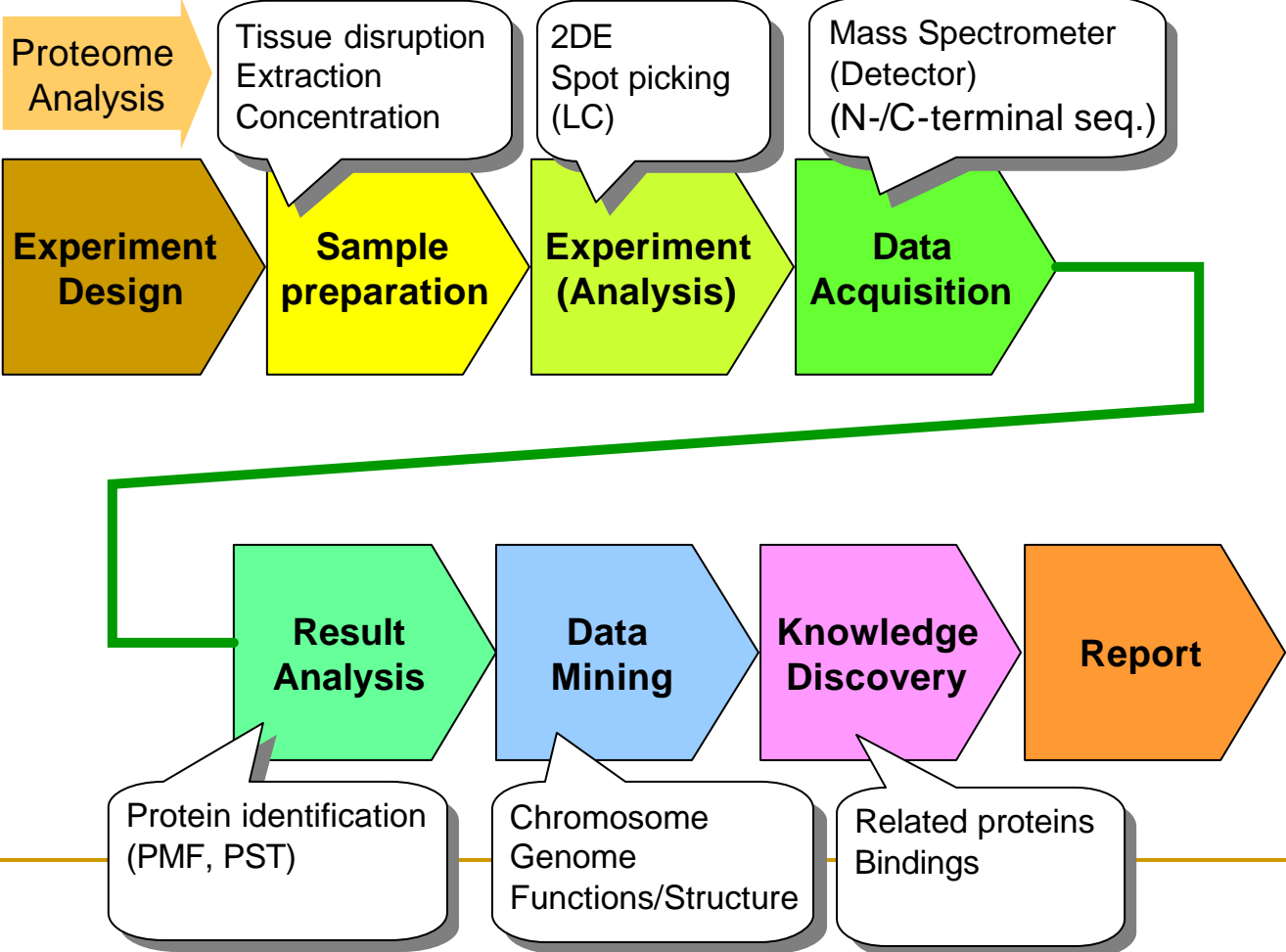Example

# XML in Bioinformatics

# Analysis flow in Life Science

Proteome Analysis

| Tissue disruption Extraction Concentration | 2DE Spot picking (LC) | Mass Spectrometer (Detector) (N-/C-terminal seq.) |

**Experiment Design** → **Sample preparation** → **Experiment (Analysis)** → **Data Acquisition**

**Result Analysis** → **Data Mining** → **Knowledge Discovery** → **Report**

Protein identification (PMF, PST)

Chromosome Genome Functions/Structure

Related proteins Bindings

# Conventional XMLs in Life Science

# Our XML-based data model

## Our XML

| Experiment Design | Sample preparation | Experiment (Analysis) |

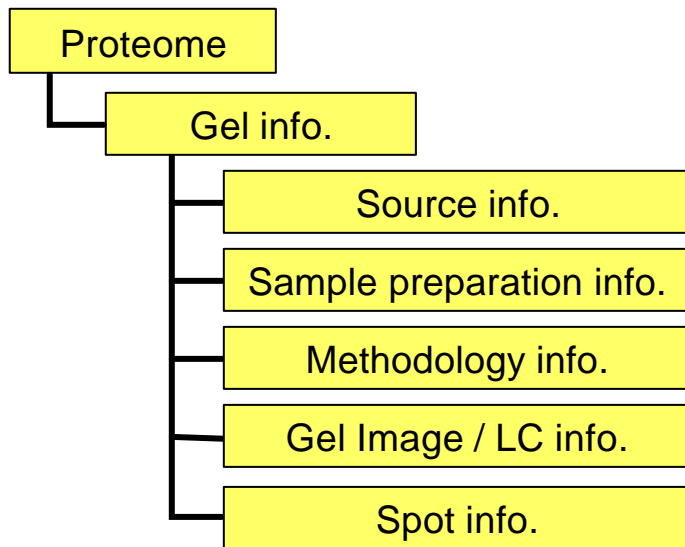| Result Analysis | Data Mining |

- **Proteome-analysis oriented**
- **Describes**
  - Sample preparation
  - Methodology
  - 2D gel image / LC results
  - Spot information
  - Sequence and feature
  - 3D structure
- **Includes other open XMLs used in life science**

**Now Available : HUP-ML (Human Proteome Markup Language) DTD and Editor**
**http://www.jhupo.org/**

# XML for Proteomics

- **Information Structure:**



```
<proteome>
  <gel   id="1">
      <source_info>
      <gel_img >
      <sample_preparation>
      <gel_conditions>
      <marker>
      <detection>
      <gel_image>
      <spot   id="1">
          ....
      <spot   id="2">

      ....
  <gel   id="2">
```
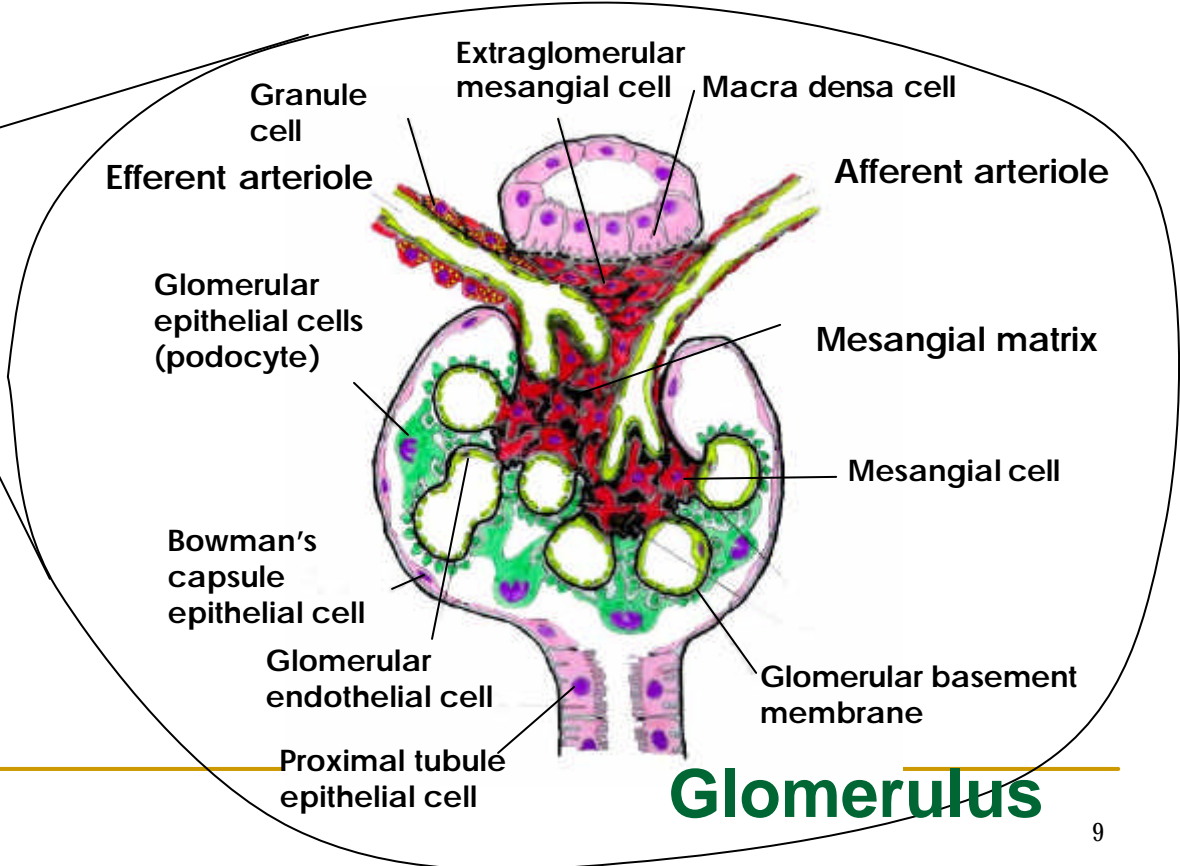
# Example:

**By A. Tsugita et al.(2002)**

## Human Kidney Glomerulus Proteome



**Nephron**

Cortex

Outer Stripe

Inner Stripe

Inner Medulla

Outer Medulla

Granule cell

Extraglomerular mesangial cell

Macra densa cell

**Efferent arteriole**

**Afferent arteriole**

Glomerular epithelial cells (podocyte)

**Mesangial matrix**

**Mesangial cell**

Bowman's capsule epithelial cell

Glomerular endothelial cell

Glomerular basement membrane

Proximal tubule epithelial cell

**Glomerulus**

# Sample of ProteomeXML (1)



**Source information**

```
- <source_info source_info_ID="HKG-1"
    creDate="2002-07-20T12:00:00"
    modDate="2002-08-10T17:20:00">
  <source>Homo sapiens</source>
  <common_name>Human</common_name>
  <strain />
  <cultiva />
  <cell_line />
  <tissue>Kidney Glomerulus</tissue>
  <plasmid />
  <growth_phase unit="year">48</growth_phase>
  <induction />
  <host />
  <description>Normal</description>
</source_info>
```

10

# Sample of ProteomeXML (2)

- <sample_preparation>
  <tissue-disruption>**Standard sieving technique using four stainless sieves. The glomeruli on the 150 micro m sieves were collected ice cold phosphate-buffered saline (PBS).**</tissue-disruption>
- <extraction>
  - <procedure>
    <process seq="**1**" action="**spin-down**"
             sample="**collection**" />
    <process seq="**2**" action="**homogenize**"
             sample="**precipitate**" >
        <add_solution solution_ID="**sol-A**"/>
    </process>
    <process seq="**3**" action="**stand**"
             time="**60**" time_unit="**min**"
             temp="**37**" temp_unit="**degree in C**" />
    <process seq="**4**" action="**centrifuge**"
             sample="**suspension**"
             time="**20**" time_unit="**min**">
        <times_g>**12000**</times_g>
    </process>

<process seq="**5**" action="**store**"
         sample="**supernatant**"
         te
         tin
  </procedure>
  <comment_ext
  </extraction>

Procedure :
(action, target, condition ) lists

- <solution solution_ID="**sol-A**" label="**2-DE lysis solution**">
  <item_solution con="**9.8**" unit="**M**" name="**Urea**" />
  <item_solution con="**2**" unit="**% w/v**" name="**NP-40**" />
  <item_solution con="**2**" unit="**% v/v**" name="**Pharmalyte(pH3-10)**" />
  <item_solution con="**10**" unit="**mM**" name="**DDT**" />
  <item_solution con="**0.5**" unit="**micro g/mL**" name="**E-64**" />
  <item_solution con="**0.5**" unit="**mM**" name="**PMSF**" />
  <item_solution con="**40**" unit="**micro g/mL**" name="**TLCK**" />
  <item_solution con="**1**" unit="**micro g/mL**" name="**aprotinin**" />
  <item_solution con="**10**" unit="**micro g/mL**" name="**chymostain**" />
  <item_solution d
  <item_solution d
  <comment_solu
  </solution>

Solution list :
solution item information

11

# Sample of ProteomeXML (3)

**Gel condition**

```
<gel_conditions gel_conditions_ID="" creDate="2002-07-20T12:00:00"
  modDate="2002-08-10T17:20:00">
- <first_dim>
 - <gel_info>
    <gel_name maker="">linear dry strip</gel_name
    <gel_pH low="3" high="10" />
    <gel_size length="24" unit="cm" />
   </gel_info>
 - <protein_solution solution_size="400" solution_unit="micro L"
     protein_amount="100" protein_unit="micro g" guiding_dye="PBP">
    <description>including standard proteins</description>
   </protein_solution>
   <rehydrate temp="20" temp_unit="degree in C" time="12" unit="hour" />
 - <running>
    <apply step="1" current="50" current_unit="micro A"
     voltage="500" voltage_unit="V" temp="20" temp_unit="degree in C"
     time="1" unit="hour" />
    <apply step="2" current="50" current_unit="micro A"
     voltage="1000" voltage_unit="V" temp="20" temp
     time="1" unit="hour" />
    <apply step="3" current="50" current_unit="micro
     voltage="8000" voltage_unit="V" temp="20" temp
     time="10" unit="hour" />
   </running>
 <IEF pH_low="3" pH_high="10" load_direction="cathode to anode" />
```

**Gel Information :**
  **Size, pH, .....**

**Running :**
**(action, condition ) lists**

# Sample of ProteomeXML (4)



PIR data area

Spot information area

# XML Editor for Proteomics Information

Our XML Document

Gel Image

Gel Info.

Spot Info.

14

# XML Editor ( Example)

Spot list

# XML Editor ( Browsing)



XML Editor

# XML Editor ( Source Information)

Source Information

<source>
<common_name>
<strain>
<cultiva>
<cell_line>
<tissue>
<plasmid>
<induction>
<host>
<growth_phase>

It is possible to import form 'templates' or other XML documents.

# Features of our data model

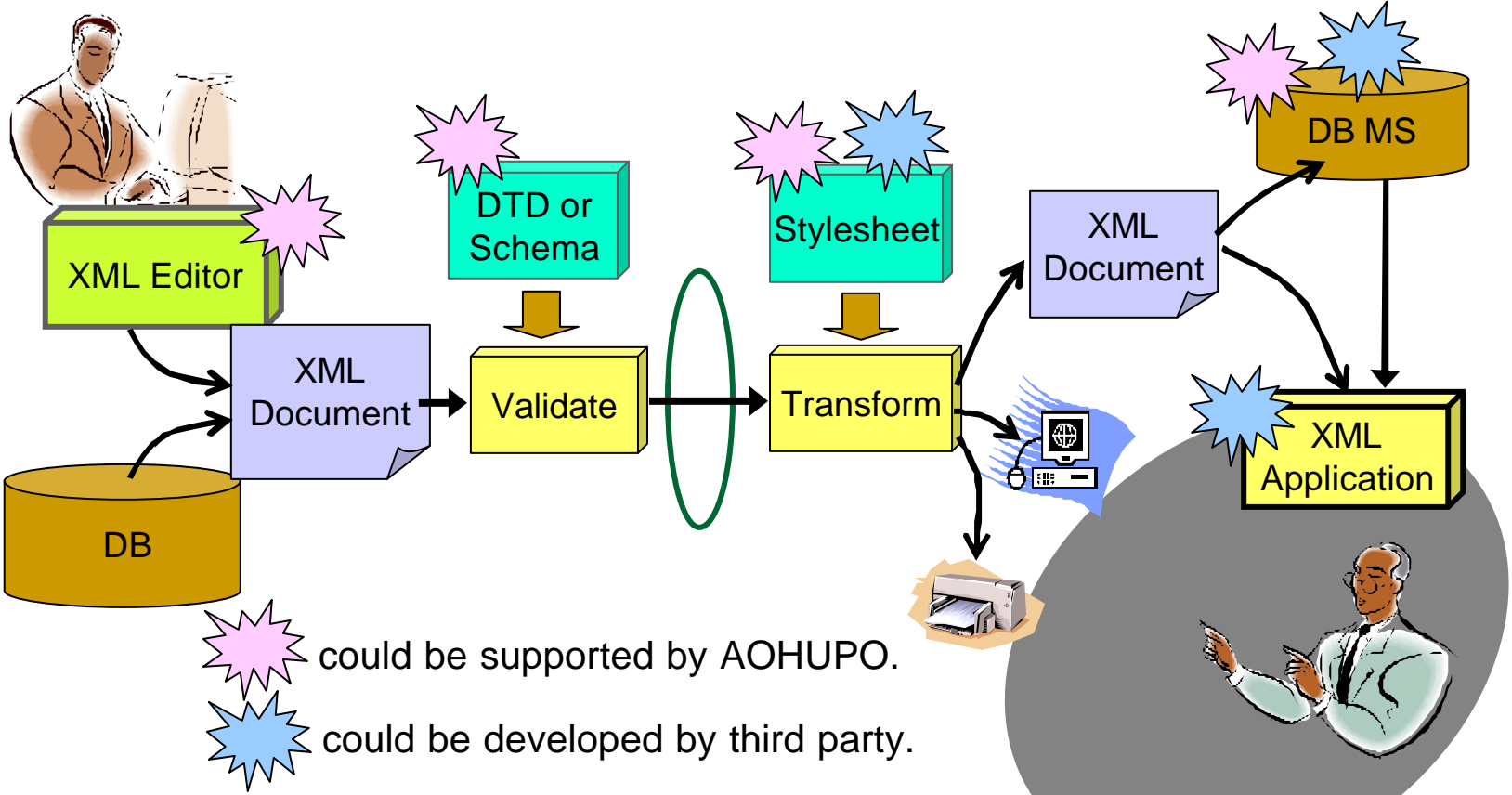Our proteomics XML:

- **describes sample preparations**
  - ❑ Improves reliability of analysis results
- **can distribute experimental information**
  - ❑ share know-how
  - ❑ improves skills
- **handle both gel-image and analysis results**
- **describes analysis information**
  - ❑ image recognition

**Now Available : HUP-ML (Human Proteome Markup Language)  DTD and Editor**
**http://www.jhupo.org/**

# Future works

- **Open DTD and/or XML Schema**
  - ❑ Collaboration with AOHUPO

- **Develop XML viewer for free distribution**

- **Prototype WWW-based management system**
  - ❑ for registration, viewing, and retrieval of entries

- **Convert from other XML formats**

- **Relation to other analysis tools**
  - ❑ image-analysis software
  - ❑ homology-analysis tools, etc.

AOHUPO: Asia Oceania Human Proteome Organiazaion

# Our XML Workflows

could be supported by AOHUPO.

could be developed by third party.

**Now Available : HUP-ML (Human Proteome Markup Language)  DTD and Editor**
**http://www.jhupo.org/**

20