# Managing large atomic and molecular data sets: HITRAN, ExoMol and CascadesDB

**IAEA**
International Atomic Energy Agency
*Atoms for Peace*

Christian Hill
Atomic and Molecular Data Unit
Nuclear Data Section
IAEA

2018 Joint ICTP-IAEA School and Workshop on
*Fundamental Methods for Atomic, Molecular and Materials Properties in Plasma Environments*

ICTP — The Abdus Salam International Centre for Theoretical Physics

# Summary

1. **Principles** of database design
2. **HITRAN** and **HITRAN*online***: low-temperature, high-resolution spectroscopic database
3. **ExoMol**: high-temperature, high-resolution spectroscopic database
4. **QuantemolDB** and **ALADDIN**: collisional databases for plasma processes
5. **CascadesDB**: collisional cascade molecular dynamics simulation database
6. **Crowdsourcing**

# Principles of database design

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

**F**indable

**A**ccessible

**I**nteroperable

**R**eusable

# Principles of database design

## Findable data

To be *findable* (meta)data must:

1. be assigned a globally-unique and persistent identifier (a URI such as a DOI)

2. registered in a searchable resource

URI = Uniform Resource Identifier

DOI = Digital Object Identifier

# Principles of database design

## Accessible data

1. To be *accessible* (meta)data must be retrievable from their identifier using a standardised communications protocol

2. the protocol (e.g. an API) must be open, free and universally implementable

3. the protocol may allow for authentication and authorisation.

# Principles of database design

## Interoperable data

To be *interoperable* (meta)data must:

1. represented in a formal, shared and broadly-applicable format

2. use vocabularies that follow FAIR principles

3. include qualified references to other (meta)data

# Principles of database design

## Interoperable data

To be *interoperable* (meta)data must:

1. represented in a formal, shared and broadly-applicable format

2. use vocabularies that follow FAIR principles

3. include qualified references to other (meta)data

Further things to consider

- **Physical Units**
- Phase conventions, reference / fiducial values
- Endianness (for binary data)
- Representation of null / missing / invalid data points

# Principles of database design

## Reusable data

To be *reusable* (meta)data must be:

1. richly described with accurate and relevant attributes

2. released with a clear data usage licence

3. associated with detailed provenance

# Principles of database design
## Authentication, Authorization, Accounting

An online database will usually implement a user-management system to:

1. Identify users (usernames, email addresses)

2. Authenticate users (login with password)

3. Account for users' activity with the database (logs)

# Principles of database design

## Practical Considerations

An online database must have

1. A stable, highly-available host server(s)
2. Software for managing users (registration, login, logout, password reset)
3. Legal terms and conditions, licence, privacy policy
4. SSL

In addition, it may have:

1. A documented API for automated access by codes, etc.
2. Contact / feedback form
3. An interface for *uploading* data

# HITRAN and HITRAN*online*

http://hitran.org

- Compilation of spectroscopic parameters for modelling radiative transfer in atmospheres

- Based at the Harvard-Smithsonian Centre for Astrophysics

- Mostly molecules, mostly at "low temperature"

- **$5 \times 10^6$** lines; 600 absorption cross sections

- 365 molecules

- 9000 registered users

# HITRAN and HITRAN*online*

# HITRAN and HITRAN*online*

http://hitran.org

- New line shape parameters (beyond Voigt)

- Many broadening species

- Pressure shifts

- Quantum numbers / labels

- Automated bibliography generation

- Uncertainties

# HITRAN and HITRAN*online*

## Relational Database Tables

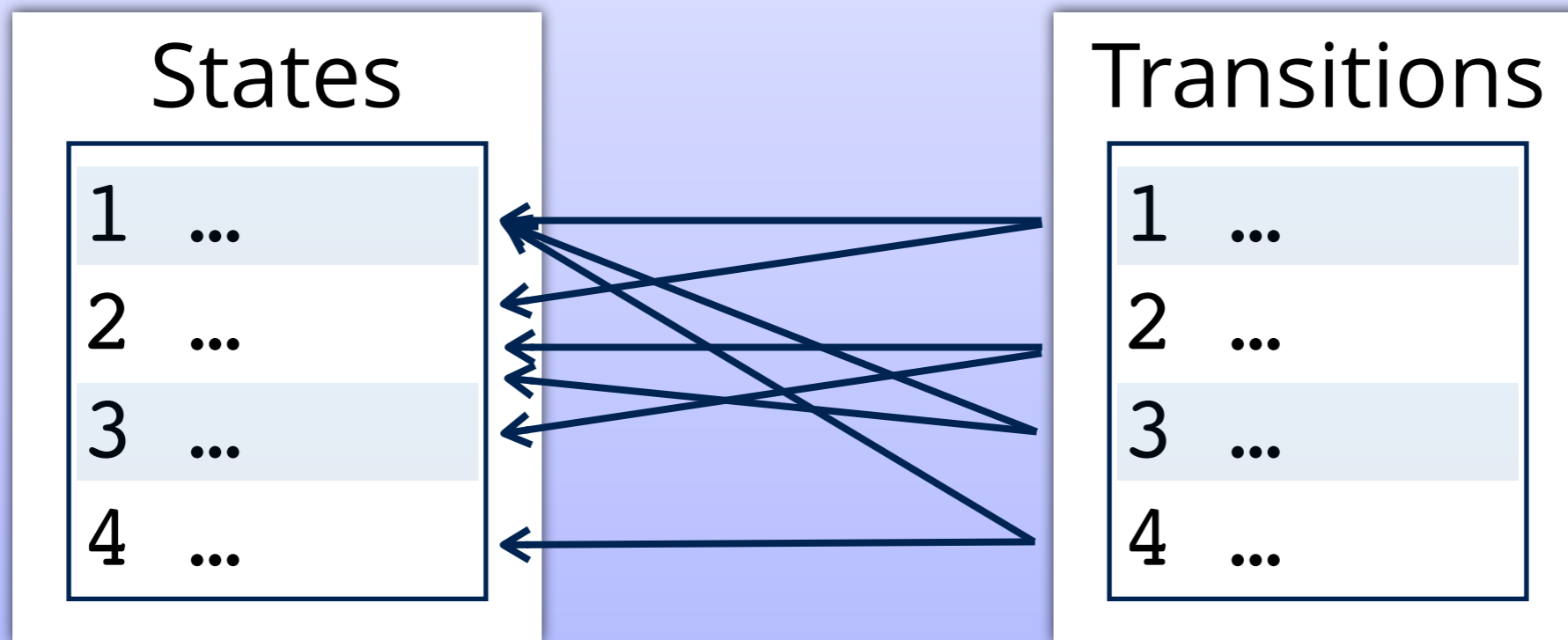Molecular **states**, linked by (radiative) **transitions**

| States | Transitions |
|---|---|
| Energy<br>Degeneracy<br>Quantum Numbers | Upper state<br>Lower state<br>*A*<br><br>... |

# HITRAN and HITRAN*online*

## User Interface

- Users register with email address and password

- Email addresses are verified

- Accessible contact form for problems / questions

- User profiles allow **customised output formats**

- Sources (citations) automatically included in output

- Interactive charts for moderate data volumes (<100,000 transitions)

# HITRAN and HITRAN*online*

## Interactive chart for data visualisation

# HITRAN and HITRAN*online*
## User-defined output formats

# HITRAN and HITRAN*online*

## HAPI

- HAPI = HITRAN Application Programming Interface
- Python-based library for accessing HITRANonline programmatically and performing common operations on the data:
  - Non-expert users can utilise advanced line shape formulations
  - Checks for updates of the latest data
  - Allows for flexible, distributable, reusable code

# HITRAN and HITRAN*online*

## HAPI

```
import hapi
hapi.db_begin('data')
hapi.fetch('CO2', 2, 1, 2000, 2100)
nu, coef = hapi.absorptionCoefficient_Lorentz(SourceTables='CO2')
plt.plot(nu, coef)
```

# ExoMol

http://exomol.com

- Compilation of molecular spectroscopic parameters for atmospheres of cool stars and exoplanets

- Some data sets can get *extremely large*:

  - $9.8 \times 10^9$ lines for $CH_4$ up to 1500 K

  - $1.68 \times 10^{10}$ lines for $PH_3$ up to 1500 K

  - $2.0 \times 10^{10}$ lines for $H_2O_2$ up to 1250 K

  - $2.1 \times 10^{10}$ lines for $SO_3$ up to 800 K

  - $6.27 \times 10^{10}$ lines for $SiH_3$ up to 1200 K

# ExoMol

## ExoMol Data Types

- *Ab initio* energy levels

- *Ab initio* transition probabilities (*A* /s$^{-1}$)

- Partition functions

- Heat capacities

- Cooling functions

- Line-by-line pressure broadening parameters

# ExoMol

## ExoMol Data Storage

Data sets too large for relational databases, so...

- Store the metadata in a relational database:

  - Molecules and Isotopologues

  - Data Types

  - Sources (citations)

- Store the energies in a single file

- Store the transitions in **compressed archives** over wavenumber intervals, with *references to the energies file* (fully normalized).

# ExoMol

## ExoMol Data Storage

e.g. $SO_3$: 2 TB → 195 GB (compressed)

states (18530508)

| | | | |
|---:|---:|---:|---:|
| 1 | 0.000000 | 1 | 0 |
| 2 | 993.679792 | 1 | 0 |
| 3 | 1059.476928 | 1 | 0 |
| 4 | 1066.497051 | 1 | 0 |
| 5 | 1591.034913 | 1 | 0 |
| 6 | 1919.634571 | 1 | 0 |
| 7 | 1981.994386 | 1 | 0 |
| 8 | 2054.050516 | 1 | 0 |
| 9 | 2061.933405 | 1 | 0 |
| 10 | 2117.465910 | 1 | 0 |
| | … | | |

transitions (21413927818)

| | | |
|---:|---:|---:|
| 10160366 | 9848857 | 2.1285e-54 |
| 10572834 | 10469949 | 1.1892e-54 |
| 1172408 | 1229247 | 5.1230e-25 |
| 1173234 | 1230094 | 4.3307e-28 |
| 12364183 | 12460001 | 5.2368e-49 |
| 12460001 | 12364183 | 5.1733e-49 |
| 1347108 | 1172690 | 1.6896e-26 |
| 150232 | 95994 | 2.8946e-30 |
| 1531681 | 1597102 | 1.7770e-35 |
| 3113447 | 3033140 | 3.6574e-54 |
| | … | |

# ExoMol

## Data Reduction

Not all applications require a full line-by-line treatment.

**Full Cross section approach**

- Pre-calculate absorption cross sections at high-resolution (wavenumber grid) for a range of *T*.

- Provide a service to interpolate and bin to the requested (*T*, Δ*v*)

$$\sigma_i = \sum_j \sigma_{ij},$$

where

$$\sigma_{ij} = \frac{S_j}{\Delta\tilde{v}} \int_{\tilde{v}_i - \Delta\tilde{v}/2}^{\tilde{v}_i + \Delta\tilde{v}/2} f_G(\tilde{v}; \tilde{v}_{0;j}, \alpha_j) d\tilde{v},$$

$$= \frac{S_j}{2\Delta\tilde{v}} \left[ \mathrm{erf}\left(x_{ij}^+\right) - \mathrm{erf}\left(x_{ij}^-\right) \right],$$

where erf is the error function and

$$x_{ij}^{\pm} = \frac{\sqrt{\ln 2}}{\alpha_j} \left[ \tilde{v}_i \pm \frac{\Delta\tilde{v}}{2} - \tilde{v}_{0;j} \right],$$

# ExoMol

## Data Reduction

**Hybrid approach**

- Use cross sections for the very large number of overlapping weak lines

- Retain line-by-line treatment for the strongest lines

S. N. Yurchenko *et al.*, *A&A* **605**, A95 (2017).

# ExoMol

## API

Per–dataset ".def" file in predefined and persistent location, e.g.

http://www.exomol.com/db/CH4/12C-1H4/YT10to10/12C-1H4__YT10to10.def

molecule

isotopologue

dataset name

Includes data version / date stamp, data column definitions, dataset file locations.

# ExoMol

## Data Expansion(!)

- Don't store redundant data:

    - Don't store more decimal places than justified by the data accuracy

    - Don't store arithmetic sequences of (e.g. wavelength or energy grids) explicitly – generate them as needed

- Don't store the data more than once:

    - Provide scripts to interconvert between commonly used formats (e.g. ExoMol → HITRAN)

# QuantemolDB

https://quantemoldb.com

Compilation of collision cross sections and rate coefficients for plasma processes:

- 16715 Reactions

- 17491 Data sets

- 1913 Species

- 3113 "Stateful Species"

QuantemolDB

Relational Database Structure (MySQL)

# QuantemolDB
## Cross section search

# QuantemolDB

## Interactive cross section comparison

# QuantemolDB

## Quantemol API

- Implemented via GET query in URL

- Authenticate users through API key

- Specify desired output format

- Return zip archive of all matching files ...

- ... or use a compatible format (COMSOL, HPEM)

- Supports queries for pre-defined "Chemistries" : validated and recommended Data Sets for particular plasma processes

# CascadesDB

## Molecular Dynamics Simulations of Collisional Cascades



Figure credit: Andrea Sand, U. Helsinki

A repository of simulations of radiation damage in materials of relevance to fusion reactor design

# CascadesDB

## Data

- Stored as `.xyz` files
- Archived into batches differing only in PKA recoil direction
- (Compressed) archive up to ~10 GB in size

# CascadesDB

## Metadata (searchable)

- Attribution
- Material parameters:
  - Lattice parameters
  - Initial crystal configuration
- Simulation details:
  - Code name and version
  - Temperature
  - Simulation time
  - Interatomic potential used

# CascadesDB

## Metadata representations



MySQL:
Searchable

XML:
Computer readable

XML Schema:
Validation

HTML: Human readable

# CascadesDB

## Metadata links to data

Metadata      URIs      Online resources

.xyz file archive

Published article

(Simulation code input file)

(Interatomic potential)

(Initial config .xyz)

# Crowdsourcing

## Types of crowdsourcing

- Creation of common goods: e.g. Wikipedia
- Carrying out micro tasks in parallel: e.g. Amazon Mechanical Turk
- Idea competitions / innovation contests
- Creative crowdsourcing: graphic design, architecture
- Crowdsolving
- Crowdfunding
- Crowdsearching
- Collaborative journalism
- Distributed computing

# Crowdsourcing

## Types of crowdsourcing

- Creation of common goods: e.g. Wikipedia
- Carrying out micro tasks in parallel: e.g. Amazon Mechanical Turk
- **Idea competitions / innovation contests**
- Creative crowdsourcing: graphic design, architecture
- Crowdsolving
- Crowdfunding
- Crowdsearching
- Collaborative journalism
- **Distributed computing**

# Crowdsourcing

## A Brief History of Crowdsourcing

1714: Find a simple and practical method for the precise determination of a ship's longitude at sea (John Harrison: H4 sea watch)

1820: First Montyon prize awarded

1884: First fascicle of the OED (800 volunteers)

1957: Jørn Utzon won the design competition for the Sydney Opera House

1999: SETI@home

2001: Wikipedia

2007: Galaxy Zoo

# Crowdsourcing

## 1. Classifying and visualising radiation damage in fusion-relevant materials

# Crowdsourcing

## 1. Classifying and visualising radiation damage in fusion-relevant materials

- Competition / Challenge with €5,000 prize
- Provided data: ~50 MD simulations (.xyz files)
  of collisional cascades:
  - Fe and W
  - different PKA energies
  - (different recoil directions)
- Scientific Leads: Andrea Sand (U. Helsinki),
- Sergei Dudarev (CCFE)
- April – June 2018

# Crowdsourcing

## 1. Classifying and visualising radiation damage in fusion-relevant materials

Participants are invited to come up with novel ways to visualise, analyse and explore the provided data. Successful submissions may involve one or more of the following:

- *Novel software for visualizing the material damage represented by the data files in a way that aids its qualitative and quantitative assessment.*

- *New software tools to rapidly and reliably identify, classify and quantify new patterns and structures of particular kinds in the data sets.*

- *Efficient algorithms to depict and summarise the statistical distribution of atom displacements and to analyse the effect of impact energy on this distribution.*

# Crowdsourcing

## Distributed Computing

- Invite members of the public to download MD simulation software to evolve a virtual crystal after impact damage.
- To be based on the BOINC (Berkeley Open Infrastructure for Network Computing) platform
- Data transferred to and from CascadesDB database

# Crowdsourcing

## Distributed Computing

**Advantages**
- Large scale, parallel computing power (cf. 600,000 users for climateprediction.net)
- Uncertainty quantification; interatomic potential validation
- Material discovery

**Challenges**
- Security
- Bandwidth, storage, scalability
- Maintaining user engagement

# Managing large atomic and molecular data sets: HITRAN, ExoMol and CascadesDB

## Thank You