# CAESAR
*(Comprehensive spAce wEather Studies for the ASPIS prototype Realization)*

The project goals are:

- Encompass all the relevant aspects of Space Weather science
- Implement a prototype for **ASPIS** (*ASI SPace Weather InfraStructure*):
  the scientific data center for Space Weather of the Italian Space Agency (**ASI**)

The project involves a large part of the **Space Weather Italian community**
bringing together 10 Italian institutions as partners:
- 3 national research institutes
  - **INAF, INGV, INFN**
- 7 universities
  - **Tor Vergata (Rome II), Perugia, Genoa, L'Aquila, Calabria, Catania, Trento**

CAESAR will aim to encompass the whole chain of phenomena from the Sun to the Earth up to planetary environments, selecting a number of well-observed "target Space Weather events" for detailed and comprehensive studies in order to showcase the proposed approach.
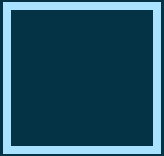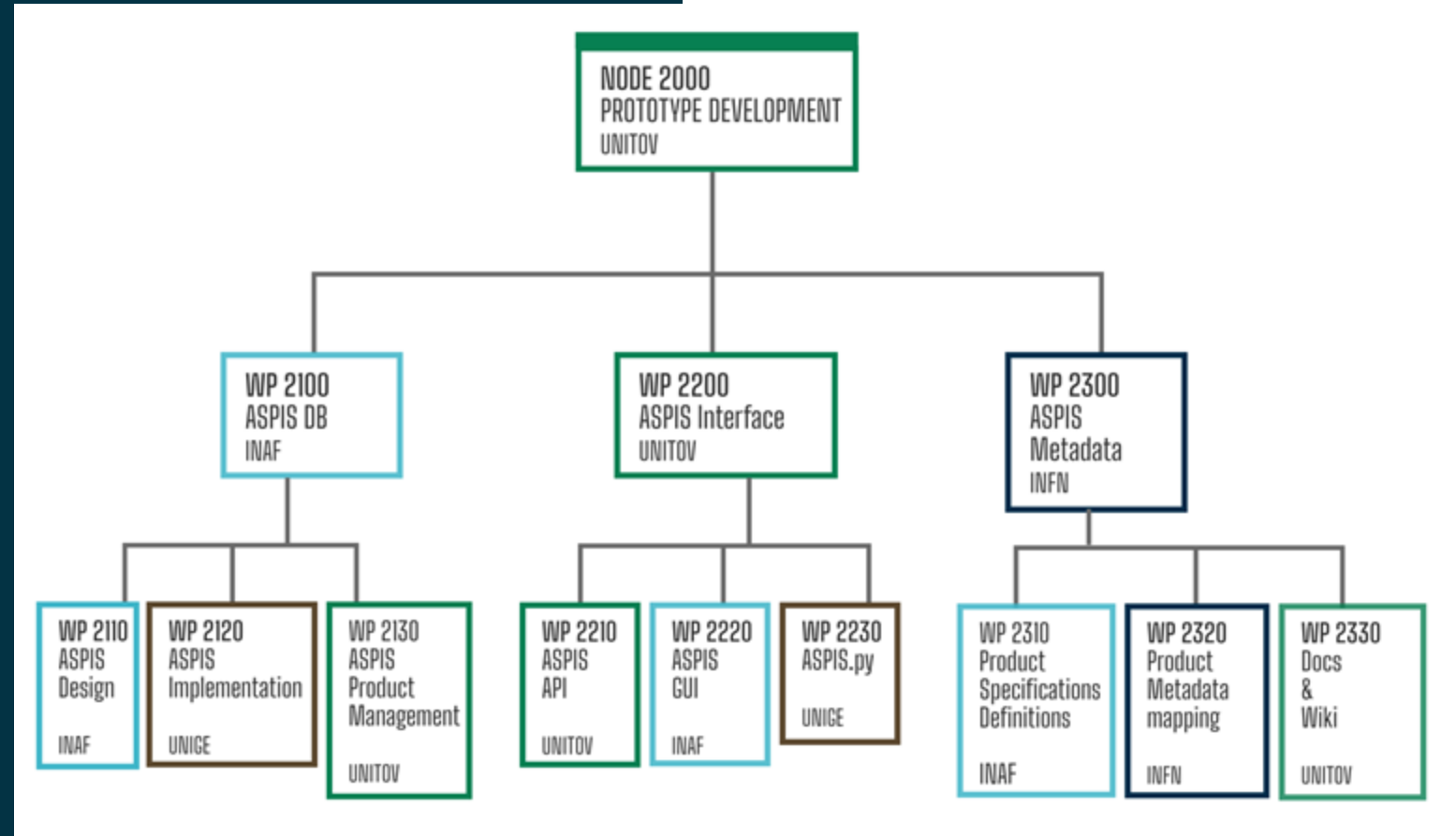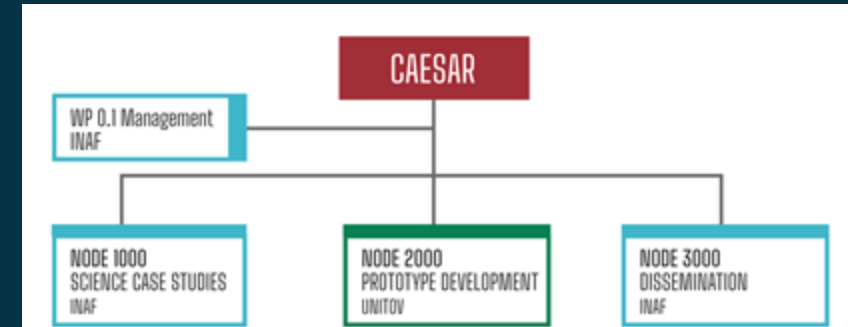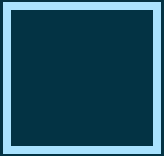
# CAESAR
*(Comprehensive spAce wEather Studies for the ASPIS prototype Realization)*



The CAESAR project is split into three main "nodes":

- Node 1000
  - <u>Science cases:</u>
    *collects all relevant use cases and data products to be included in the project*
- Node 2000
  - <u>Prototype development:</u>
    *study and design the actual implementation of the platform*
- Node 3000
  - <u>Dissemination:</u>
    *outreach activities and raise awareness on the CAESAR/ASPIS projects*
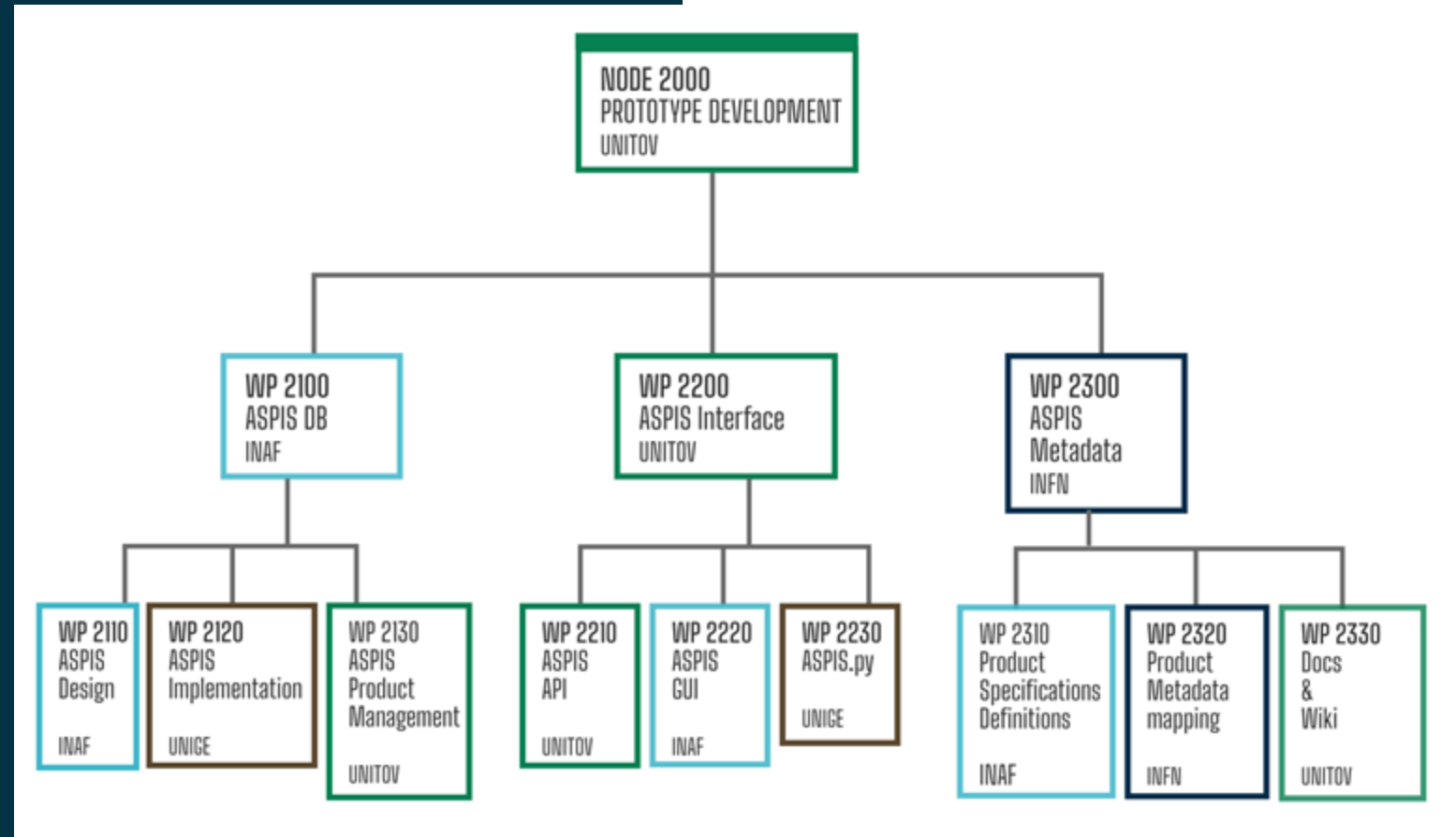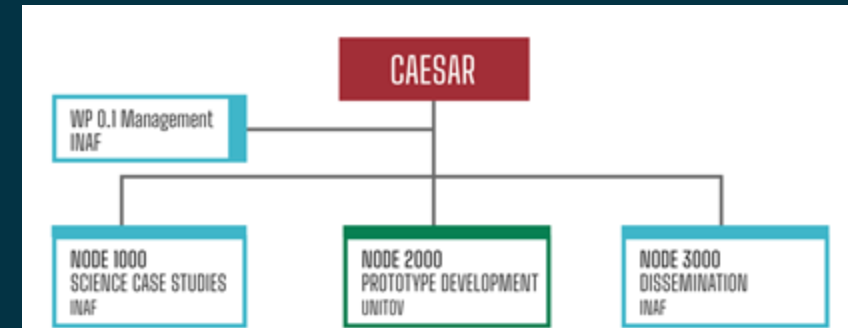
# CAESAR *(Node 2000)*

*(Comprehensive spAce wEather Studies for the ASPIS prototype Realization)*



On the technical side the prototype implementation will be based on three main pillars:

- Database [WP2100]
- Interface [WP2200]
- Metadata mapping [WP2300]

The ASPIS DB will contain mainly proprietary/co-proprietary products, with their relative data policy in a homogeneous, standardized collection of resources.

Other important external data shall be accessed through links to existent archives.
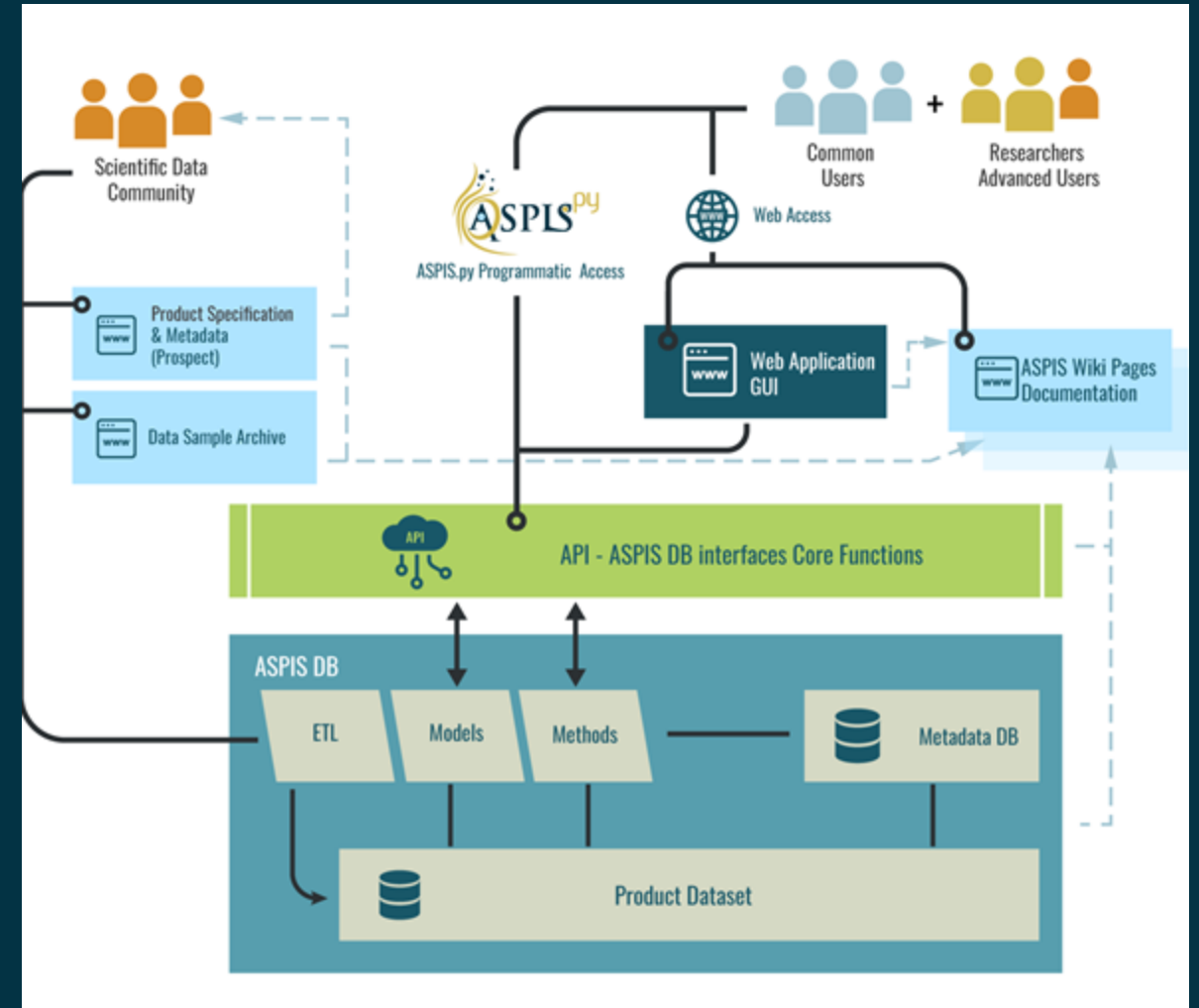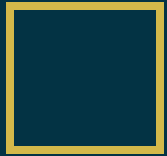
# CAESAR *(Node 2000)*
*(Comprehensive spAce wEather Studies for the ASPIS prototype Realization)*

High Level Working Packages are responsible for
- WP2100
  - Designing the DB
    *(database in its general meaning)*
  - Implement it
    - Empty → preliminary → final prototype
  - Ingest product data/models
  - Incorporate models/tools
- WP2200
  - Define the archive's API
  - Provide a GUI (available to all user)
  - Develop ASPIS.py (for advanced researchers)
- WP2300
  - Templating product descriptions
    *(and collect them from Node 1000)*
  - Map metadata content and formats
    *(for internal/external usage)*
  - Document all the activities and processes

# Choosing a data format

Before collecting product specifications from users in Node 1000 we needed to choose a unique data format to store them

Several factors were taken into account:
- We have a high number of products, coming from several physics domains and each with its own custom data format
- It is very unlikely to have common fields or variable names between different products
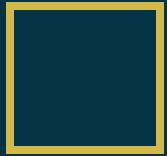
# Choosing a data format

Before collecting product specifications from users in Node 1000 we needed to choose a unique data format to store them

Several factors were taken into account:
- We have a high number of products, coming from several physics domains and each with its own custom data format
- It is very unlikely to have common fields or variable names between different products
- Data structures can be arbitrarily nested or complex, with variable size or content

# Choosing a data format

Before collecting product specifications from users in Node 1000 we needed to choose a unique data format to store them

Several factors were taken into account:
- We have a high number of products, coming from several physics domains and each with its own custom data format
- It is very unlikely to have common fields or variable names between different products
- Data structures can be arbitrarily nested or complex, with variable size or content

<u>This pretty much excludes a tabular data format</u>

# Choosing a data format

The choice fell on the **JSON** (JavaScript Object Notation) format

It is represented as a collection of key-value pairs, where the value can store several data types, including other JSON documents (which allows arbitrary nesting of data)

It's both very flexible and powerful, and it allows to store data as strings, numbers, sequence of other values (arrays), or other JSON documents.
It was designed to be human-readable, since each value in the document is effectively "named" by its key.

```json
{
  "firstName": "John",
  "lastName": "Smith",
  "isAlive": true,
  "age": 27,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "office",
      "number": "646 555-4567"
    }
  ],
  "children": [
    "Catherine",
    "Thomas",
    "Trevor"
  ],
  "spouse": null
}
```

# Metadata Design

With a data format chosen, we then proceeded to lay down a set of crucial metadata that would help catalogue and efficiently search among products including:

- A title, shortname and unique identifier, and a type determining if the product is a dataset or a scientific model
- A "status" indicating if the product is actively maintained and updated
- Info about the product "curator" institution, which includes names of contributors and a contact person
- An accurate and concise description of the product data and its layout
- A complete set of relationships with other products
- …and much more

```
{
  "product": {
    "title": "TSST H-alpha Full Disk Images",
    "shortname": "TSST-halpha",
    "type": "data",
    "identifier": "aspis:///unitov/halpha",
    "altidentifier": [
      "doi: \thttps://doi.org/10.1051/swsc/2020061"
    ],
    "status": "Active",
    "created": "2022-04-29T12:32:27+02:00",
    "updated": "2022-04-29T12:32:46+02:00"
  },
  "template": {
    "version": 1
  },
  "curation": {
    "publisher": "University of Rome Tor Vergata",
    "publisherID": "aspis://unitov",
    "creator": [
      "Luca Giovannelli"
    ],
    "contributor": [
      "Francesco Berrilli",
      "Dario del Moro"
    ],
    "contact": [
      {
        "name": "Luca Giovannelli",
        "email": "luca.giovannelli@roma2.infn.it",
        "address": "Via della Ricerca Scientifica 1, 00133, Roma",
        "telephone": "0672594552"
      }
    ]
  },
  "content": {
    "description": "The Ha channel of the TSST is a Daystar SR-127 QT, an achromatic refractor with
    "referenceURL": "https://helio.roma2.infn.it/",
    "relationship": [
      {
        "type": "related-to",
        "relatedproductid": "aspis:///unitov/vel",
        "relatedproduct": "TSST-vel",
        "description": "TSST Full Disk Velocity Maps"
      },
      {
        "type": "related-to",
        "description": "TSST Full Disk LoS magnetic maps",
        "relatedproduct": "TSST-mag",
        "relatedproductid": "aspis:///unitov/mag"
      },
      {
        "type": "related-to",
        "description": "TSST Visualization codes",
        "relatedproduct": "TSST-vis-codes",
        "relatedproductid": "aspis:///unitov/vis-codes"
      },
      {
        "type": "related-to",
```

# Pros / Cons

JSON files are easy to inspect by humans (useful for Node 2000 people during the prototype development)

JSON files are easy to handle by computers, allowing Node 2000 to easily automate part of their tasks (e.g. generation of wiki entry for each product, DB schema design, etc...)



JSON can be difficult to write (and time consuming) especially for first-time users, which includes most if not all the data providers in Node 1000

In order to make this process as streamlined and intuitive as possible we developed a dedicated web application that Node 1000 users could use to fill in their product metadata and save them directly in JSON format

# ProSpecT

The **ProSpecT** (Product Specification Template) project is based on the JSONForms web framework (https://jsonforms.io)

It allows users to input all the product metadata in pre-determined fields and provides some basic form of input validation to make sure the metadata are in the right form

Metadata are divided in sections to ease navigation and to not overwhelm the user with dozens of fields on the same page
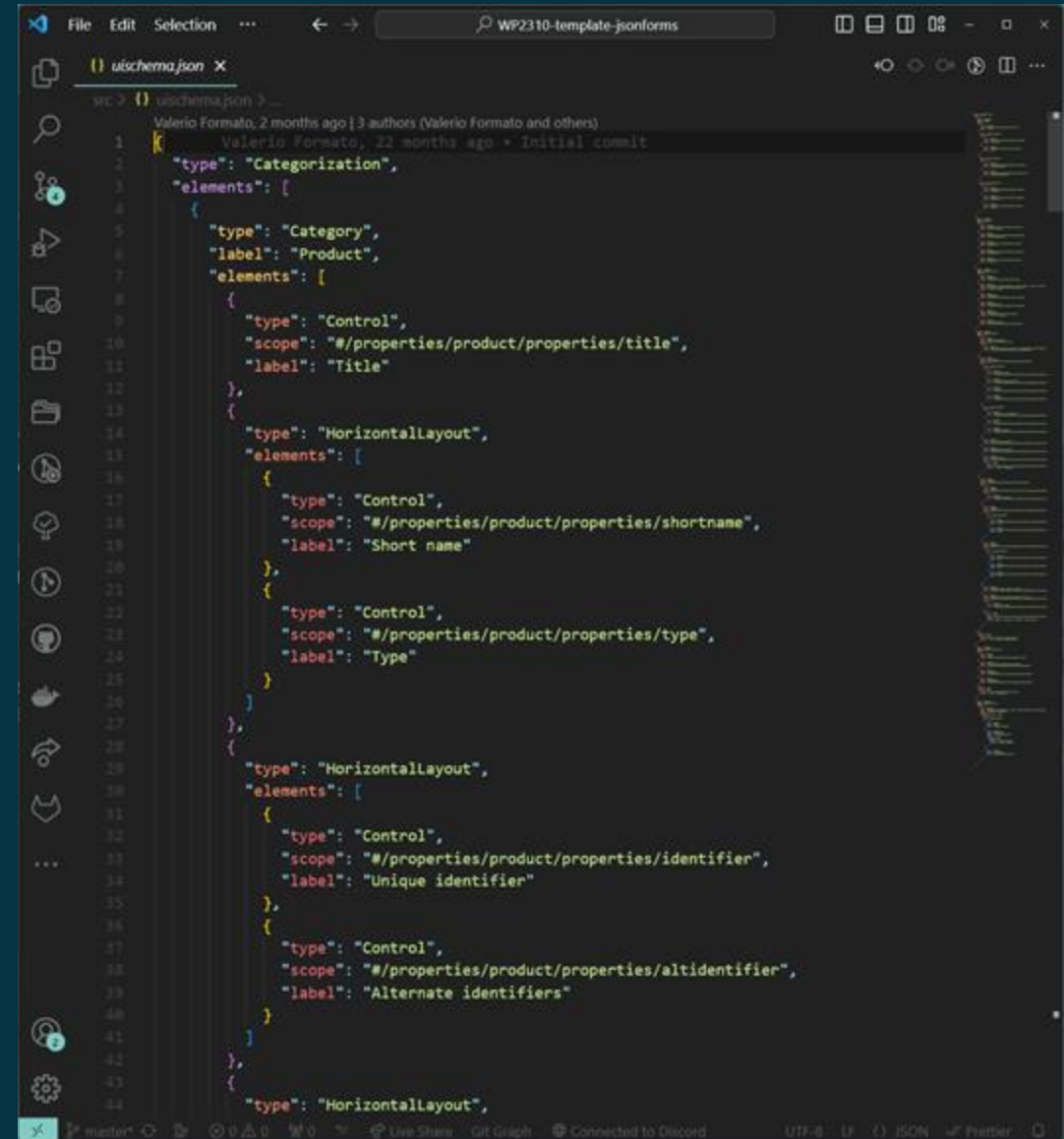
# ProSpecT

The data in the forms will be stored in a JSON object.

JSONForms will use a special "schema.json" file to infer the final layout of the output JSON forms content.

# ProSpecT

The data in the forms will be stored in a JSON object.

JSONForms will use a special "schema.json" file to infer the final layout of the output JSON forms content.

Similarly, JSONForms requires a "uischema.json" file which specifies where and how the UI elements should be rendered in the resulting webpage.

UI elements are linked to a specific field in the JSON output via the special "scope" key.

# ProSpecT

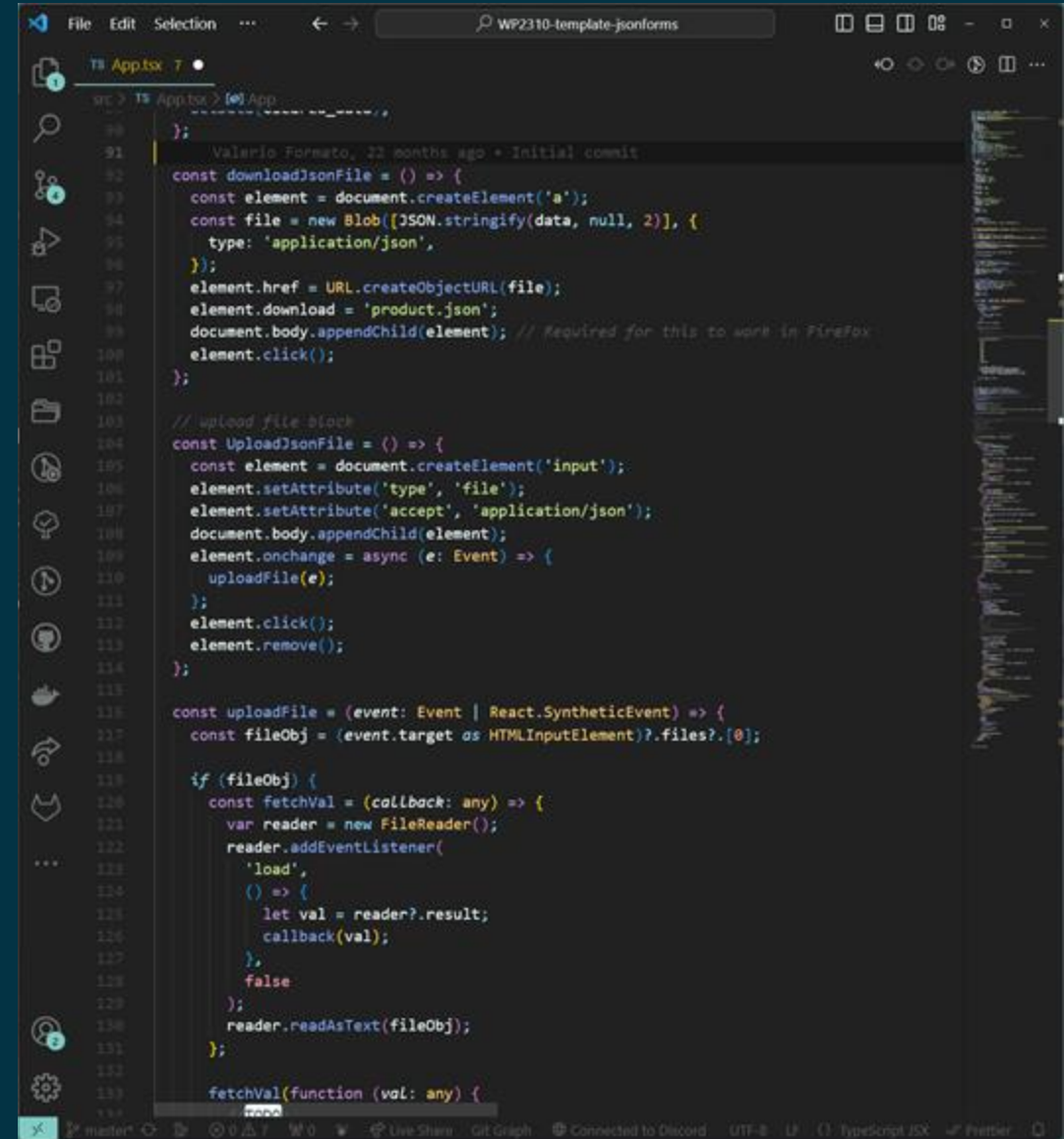The data in the forms will be stored in a JSON object.

JSONForms will use a special "schema.json" file to infer the final layout of the output JSON forms content.

Similarly, JSONForms requires a "uischema.json" file which specifies where and how the UI elements should be rendered in the resulting webpage.

UI elements are linked to a specific field in the JSON output via the special "scope" key.

JSONForms supports different JS frameworks and custom features can be added to the main application.
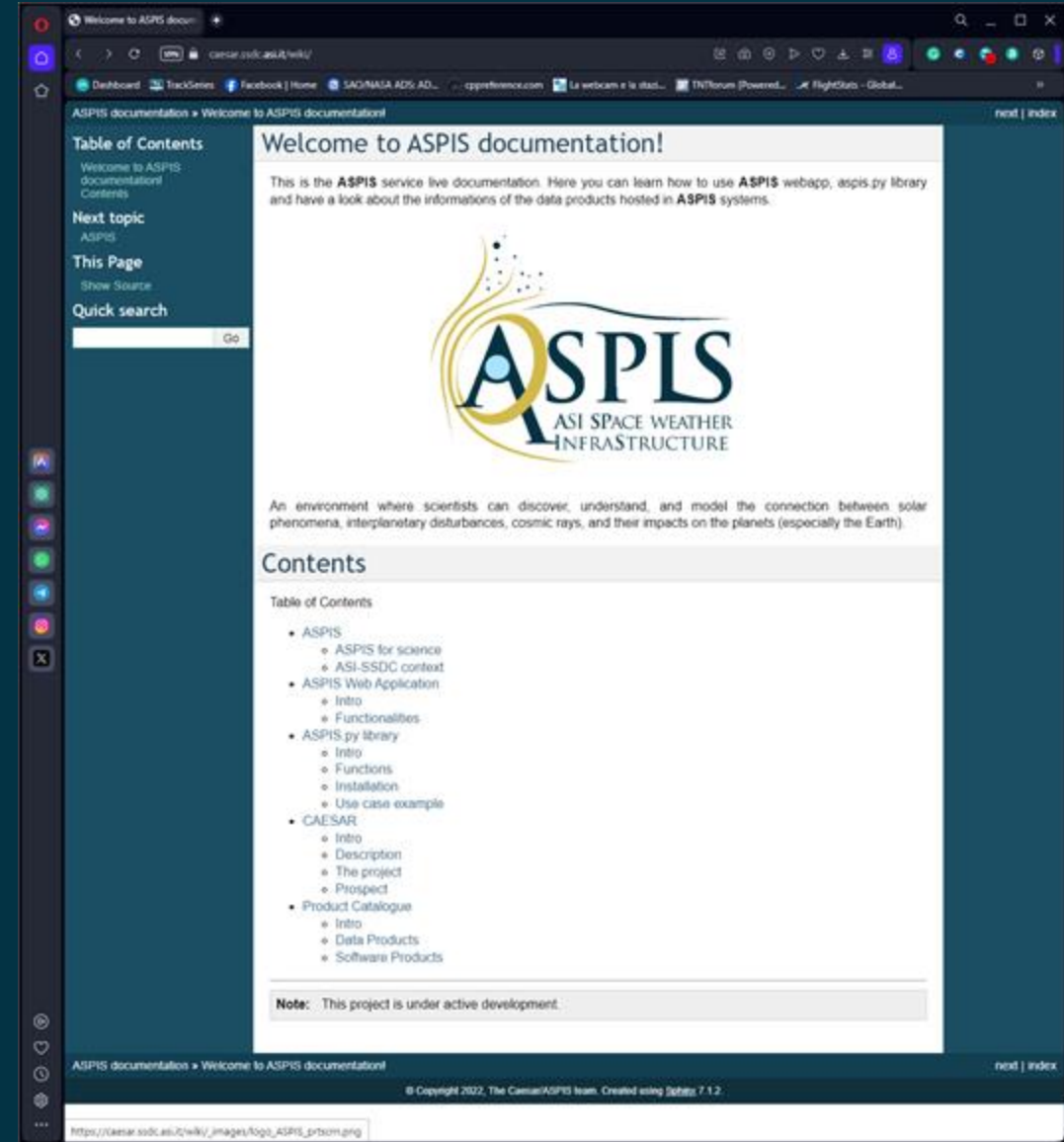*e.g. ProSpecT users can save an incomplete form, and re-upload it later to finish working on it.*

# ProSpecT

One of the results of this effort can be see in the accompanying wiki, which also holds a list of all the products included in the ASPIS prototype database
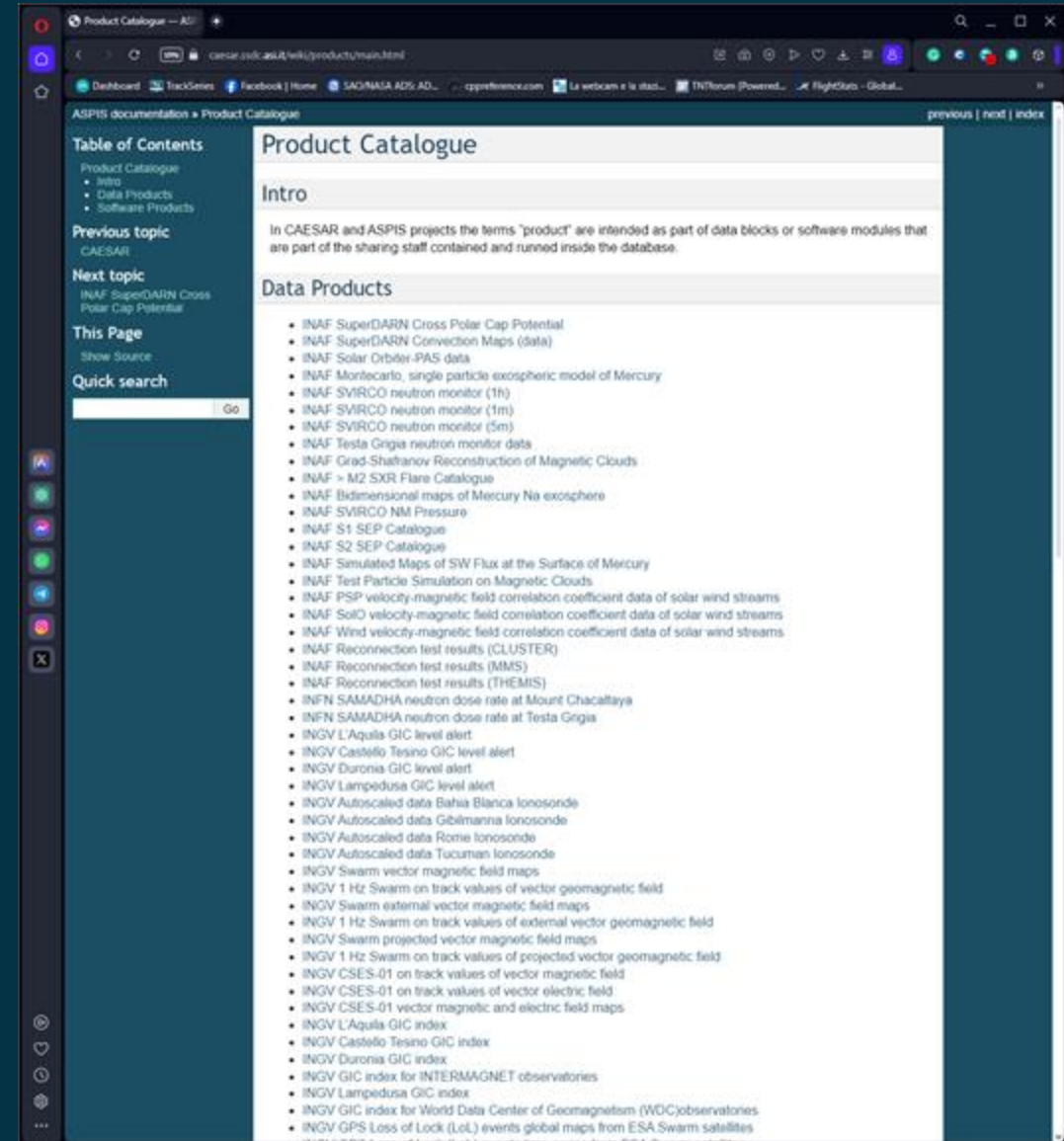
# ProSpecT

One of the results of this effort can be see in the accompanying wiki, which also holds a list of all the products included in the ASPIS prototype database

The full list is generated from the collection of JSON metadata files that data providers in Node 1000 created using ProSpecT

# Conclusions

The CAESAR project was created with the specific goal of implementing a prototype for the ASPIS project.

A large number of science cases and data products will be handled by the project. To facilitate the design process, and in future help querying for data products in ASPIS, a robust set of metadata should accompany each product.

JSON was chosen as the data format for metadata description.
To facilitate creation of such metadata files a dedicated tool (ProSpecT) was created using the JSONForms framework.

The resulting metadata files have been successfully used to automate large tasks (such as the creation of documentation) and are still largely used in the final design of the ASPIS prototype database.

23/11/2023  ESWW 2023

# Thank you!

**CAESAR NODE 2000 Team**
Valerio Formato (INFN), Marco Molinaro (INAF), Valeria di Felice (INFN), Dario Del Moro (UNITOV), Stefano Scardigli (UNITOV), Cristina Campi (UNIGE), Federico Benvenuto (UNIGE), Carmelo Magnafico (INAF), Rossana De Marco (INAF), Ermanno Pietropaolo (UNIAQ), Gregoire Francisco (UNITOV), Andrea Tacchino (UNIGE), Mirko Stumpo (INAF), Liu Scigè John (INAF), Giuseppe Di persio (INAF), Emanuele Scalise (INAF), Raffaella Noschese (INAF), Monica Laurenza (INAF)

**EB Team**
Cristina Plainaki, Anna Milillo, Giuseppe Sindoni, Marco Giardino, Alberto Bigazzi, Gianluca Polenta

Project Prime:

Project Partners: