Sveriges lantbruksuniversitet
Swedish University of Agricultural Sciences
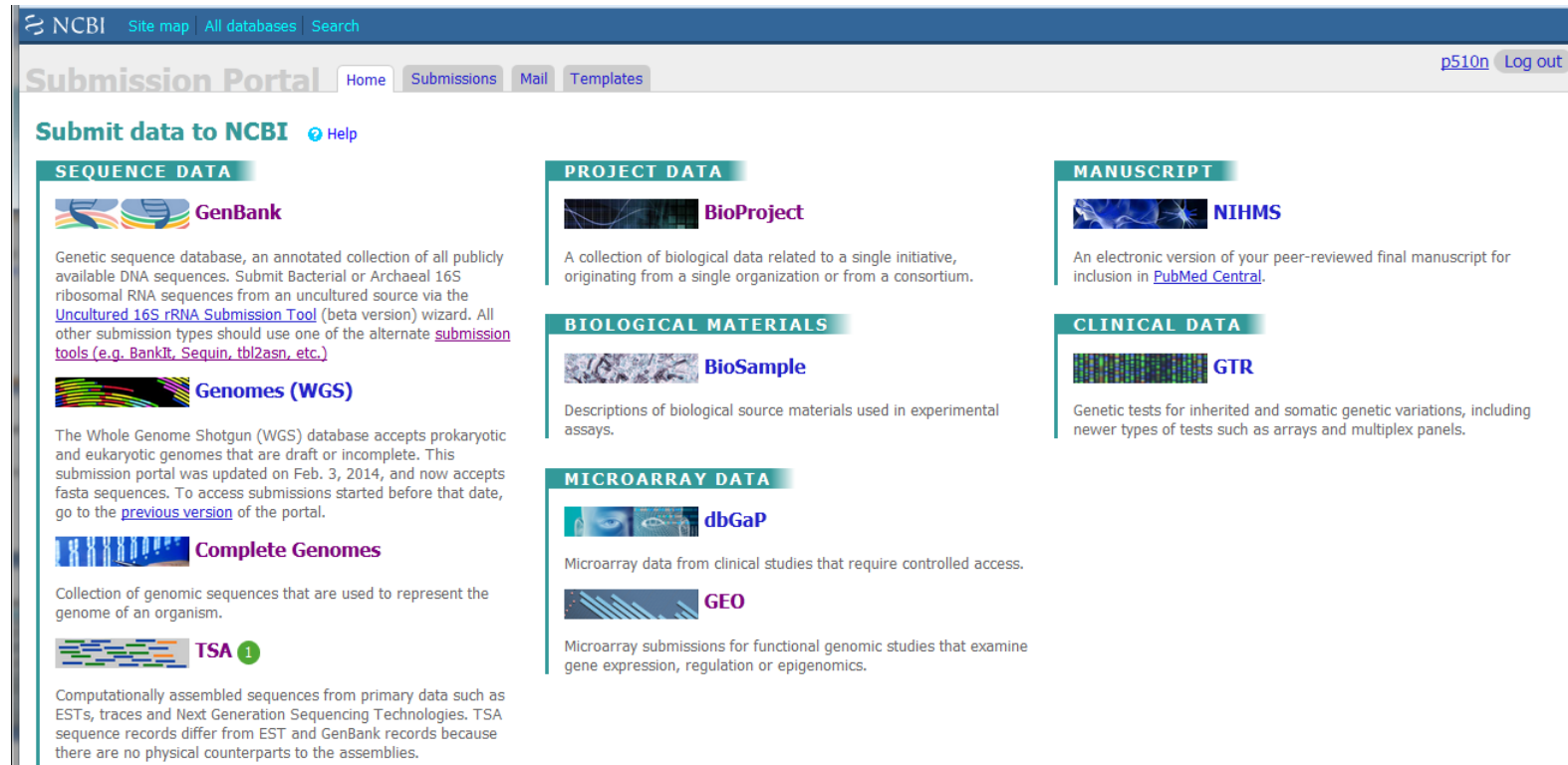
**Isak Sylvin**

# Data submission

Repositories, tools and types

# Where do I start?

https://submit.ncbi.nlm.nih.gov/

# Main data categories

- WGS and Genomes
- "Annotated sequences"
    - DDBJ/EMBL/Genbank
        - Continuous
        - Nucleotide data
        - >200 nucleotides
        - Physical counterpart
- Microarray
    - GEO
- SRA
- dbSNP & dbVar
    - <50 and >50 BP seq respectively
- Data for specific purposes (Clinical etc)

# Main metadata categories

- BioSample
    - Biological source materials
    - Enviroment

- Bioproject
    - Details about the lab samples
    - Includes BioSample

# Genbank submission tools

| Tool | Function |
|------|----------|
| Bankit | Web-based, easy, weak |
| Sequin | Stand-alone, graphical, easy, works offline |
| Tbl2asn | Stand-alone, command line, strong, offline; **Good for big data** |
| Barcode | Like Bankit but for [Barcode of Life](#) projects based on the COI gene. |

# About project size

- There are always batch alternatives
- Only use graphical tools on small projects
- Or when unsure

# About the hands-on

You will test uploading to:
- Genbank, BankIt
- Genbank, Sequin
- BioSample
- BioProject
- GEO
- SRA

# About the hands-on

1. We will not submit anything
2. Play around with as many fields as possible
3. Data can be uploaded to other sections
4. If you are unsure about what a field does; **ASK!**