

F.A.I.R. DATA MANAGEMENT WITH DATALAD

1

NIKI RUNS A TEAM OF DATA SCIENTISTS AND ENGINEERS AT A RESEARCH AND DEVELOPMENT GROUP.



2

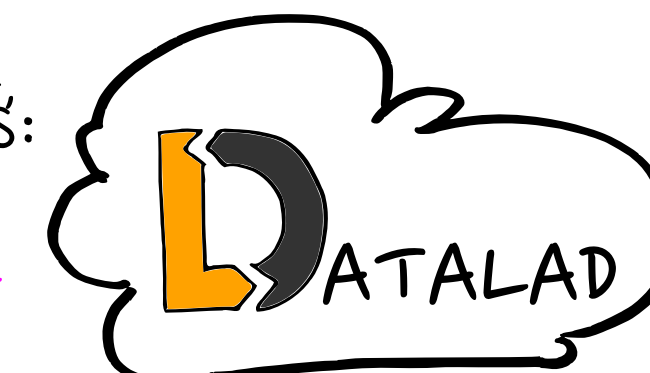


THEY WANT TO ORGANIZE AND KEEP TRACK OF ALL THEIR DATA ON A LOCAL COMPUTE CLUSTER, COMPUTERS, AND HARDDRIVES. LOW MAINTENANCE AND LOW COSTS ARE IMPORTANT.



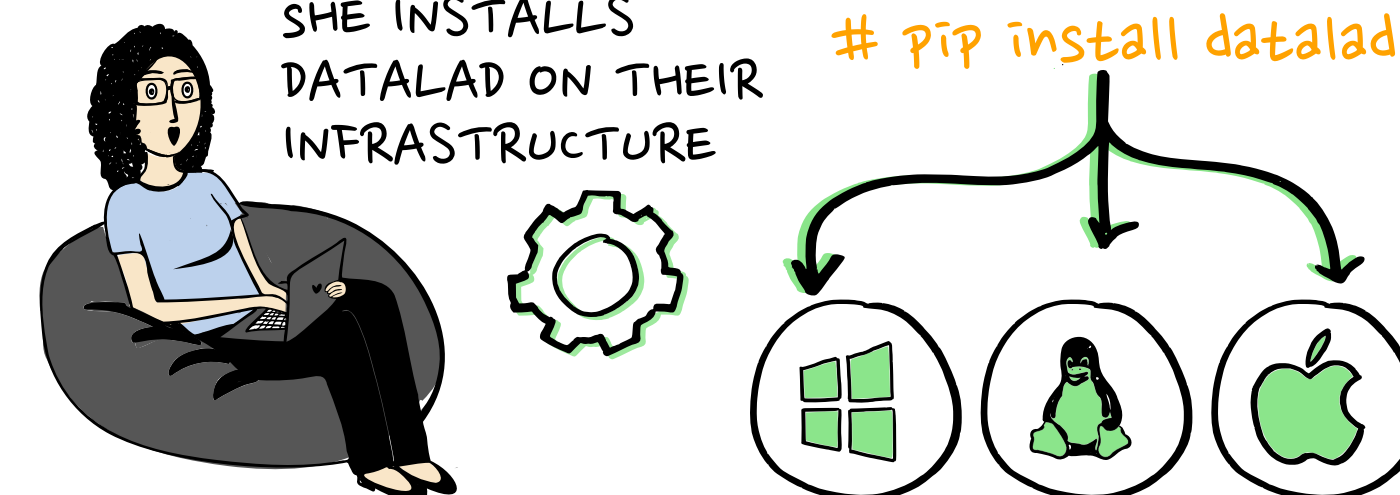
3

NIKI FINDS DATALAD ONLINE NOTICING THAT IT PROVIDES:
✓ FREE AND OPEN SOURCE
✓ DATA VERSION CONTROL
✓ & PROVENANCE TRACKING



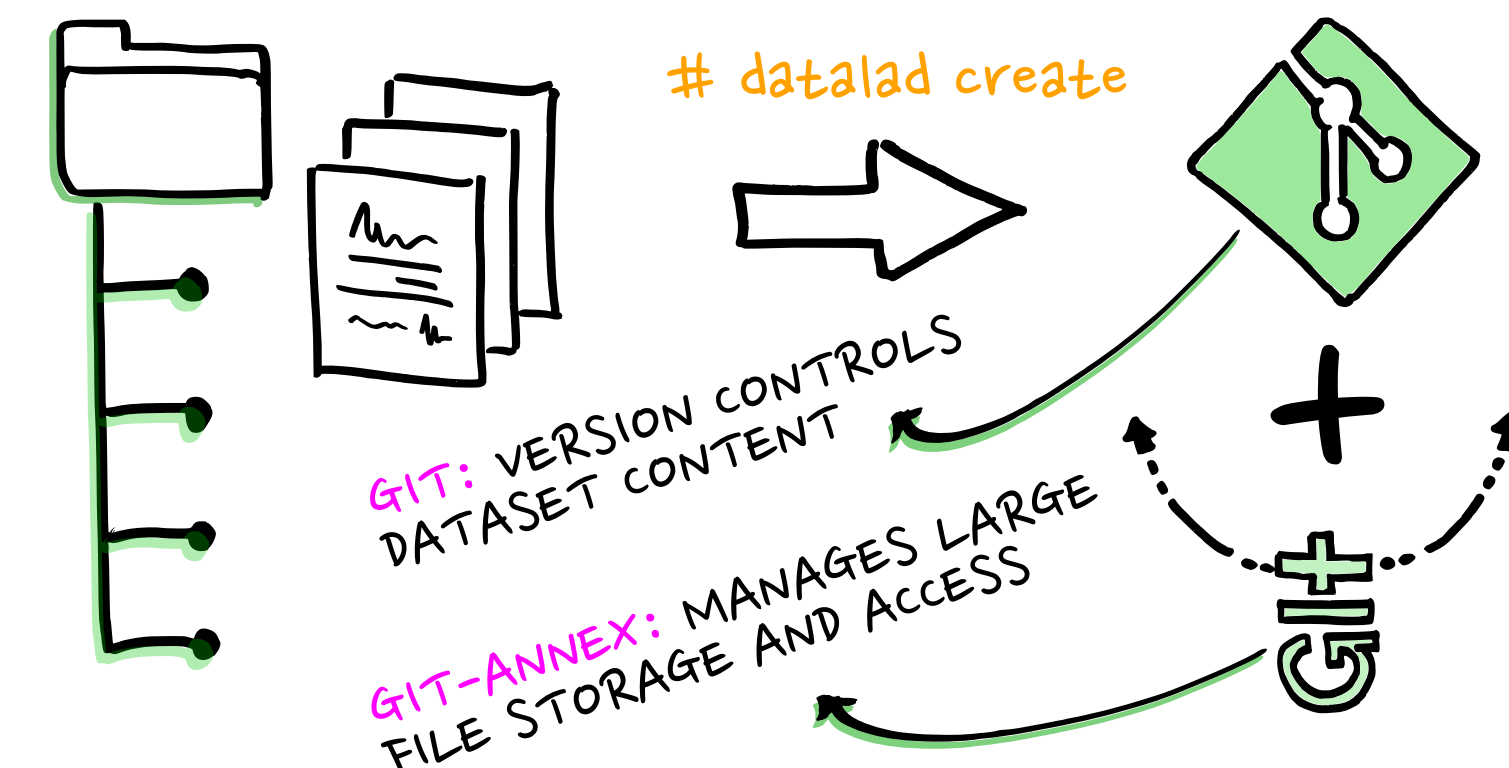
SHE INSTALLS DATALAD ON THEIR INFRASTRUCTURE

pip install datalad

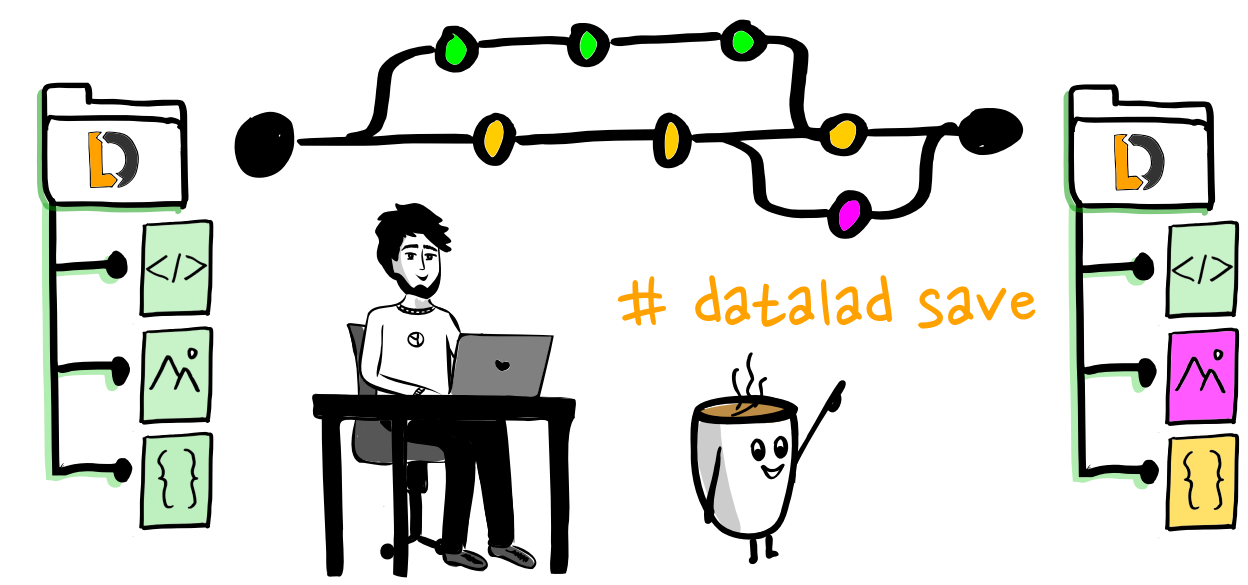


4

FIRST THEY ORGANIZE FILES INTO MODULAR UNITS AND TURN THESE INTO DATALAD DATASETS:



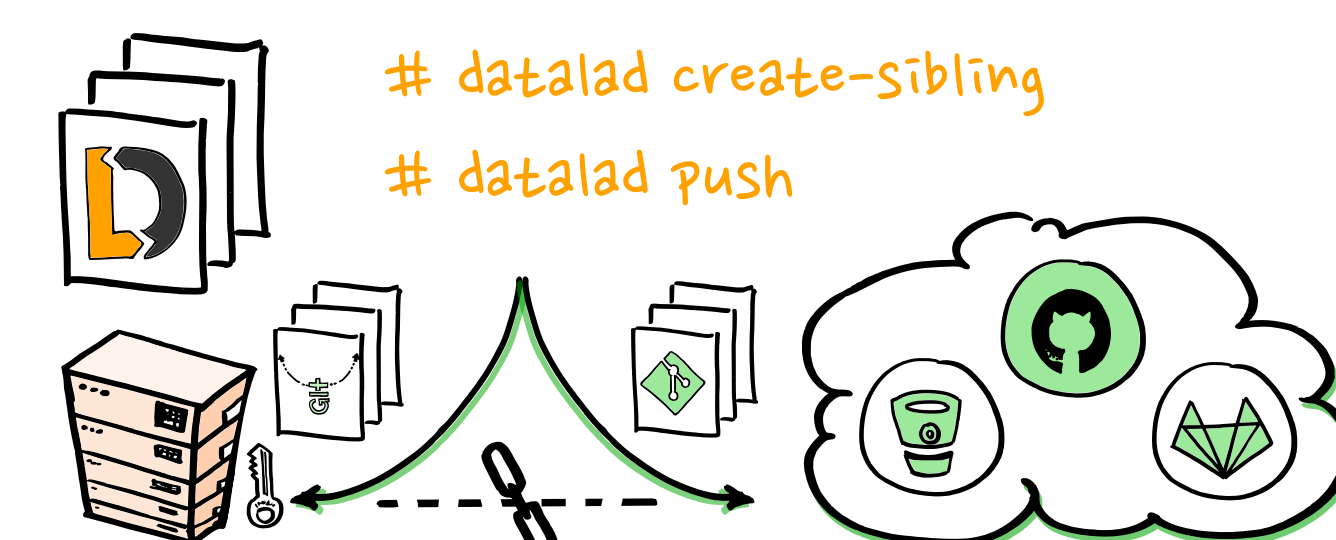
5



THEN THE TEAM LEARNS TO COLLABORATE USING DATALAD & GIT. THEY CREATE DATASET BRANCHES AND ADD, SAVE, AND MERGE CHANGES, WHILE RETAINING A FULL, AUDITABLE HISTORY.

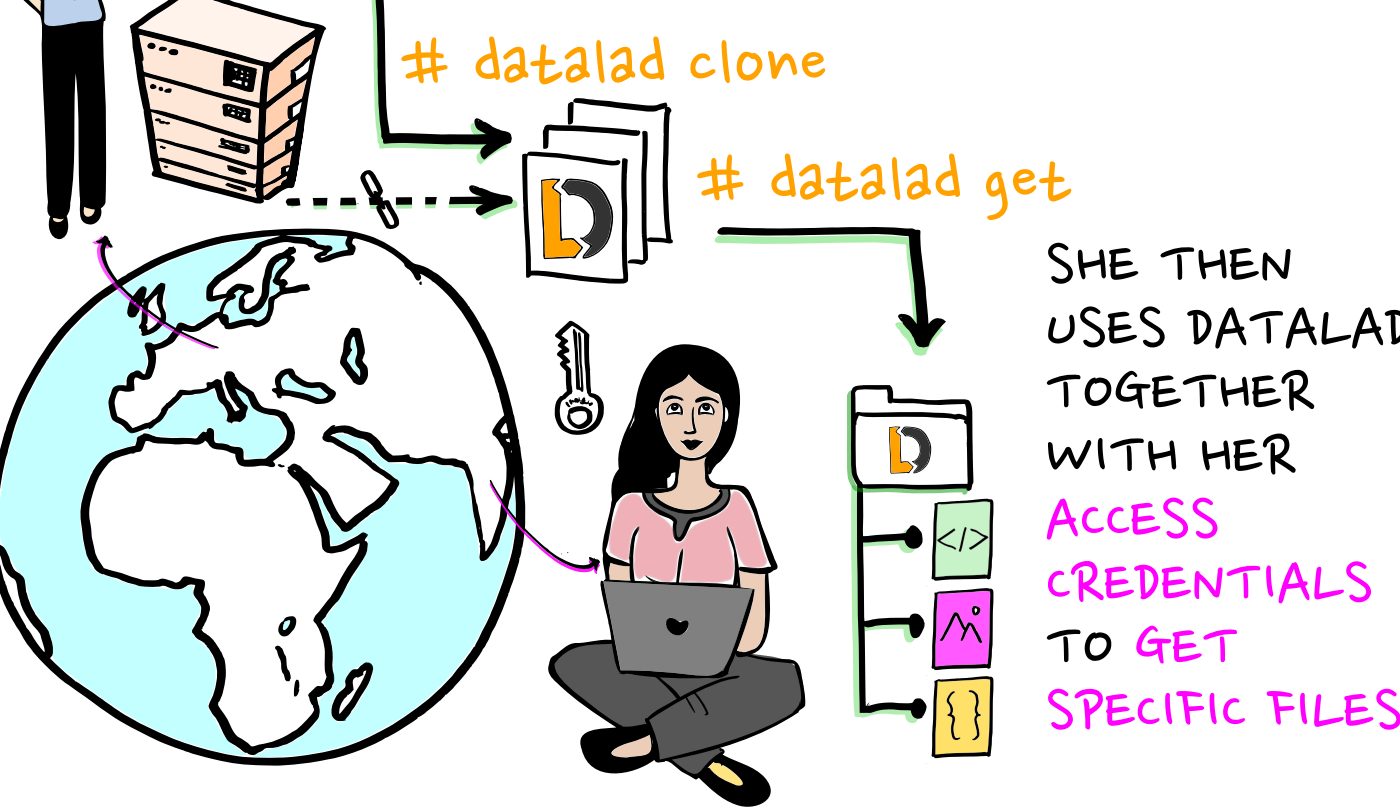
6

FOR DECENTRALIZED COLLABORATIONS, NIKI PUBLISHES THEIR DATALAD DATASETS TO GITHUB. THE LIGHTWEIGHT GIT REPOSITORIES ARE MADE PUBLIC, WHILE THE ACTUAL DATA REMAIN STORED LOCALLY, SECURELY, YET STILL ACCESSIBLY.



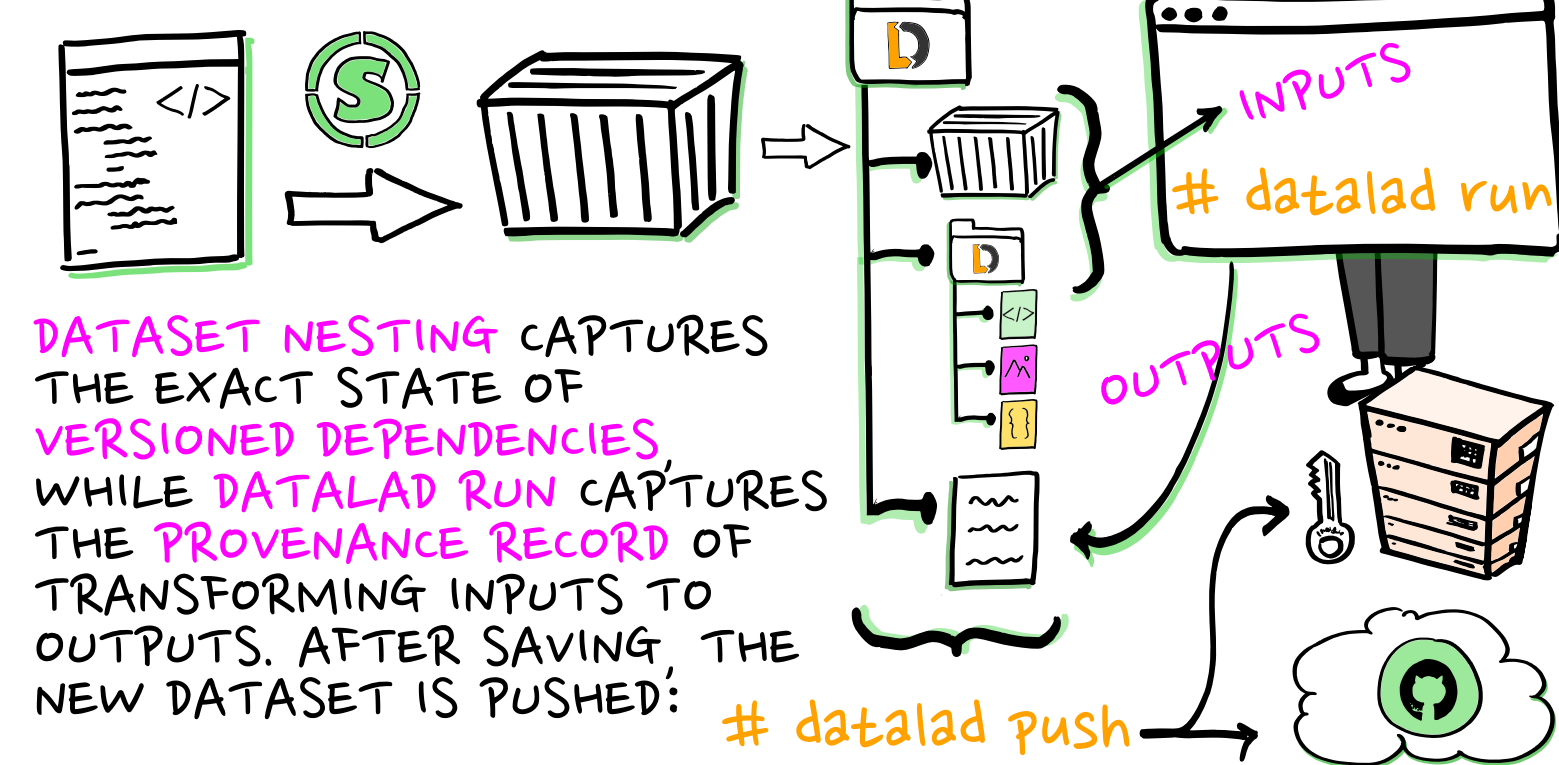
7

NIKI'S COLLEAGUE, PRIYA, CAN CLONE THE DATASET FROM GITHUB TO GET LOCAL ACCESS TO THE FILE TREE.

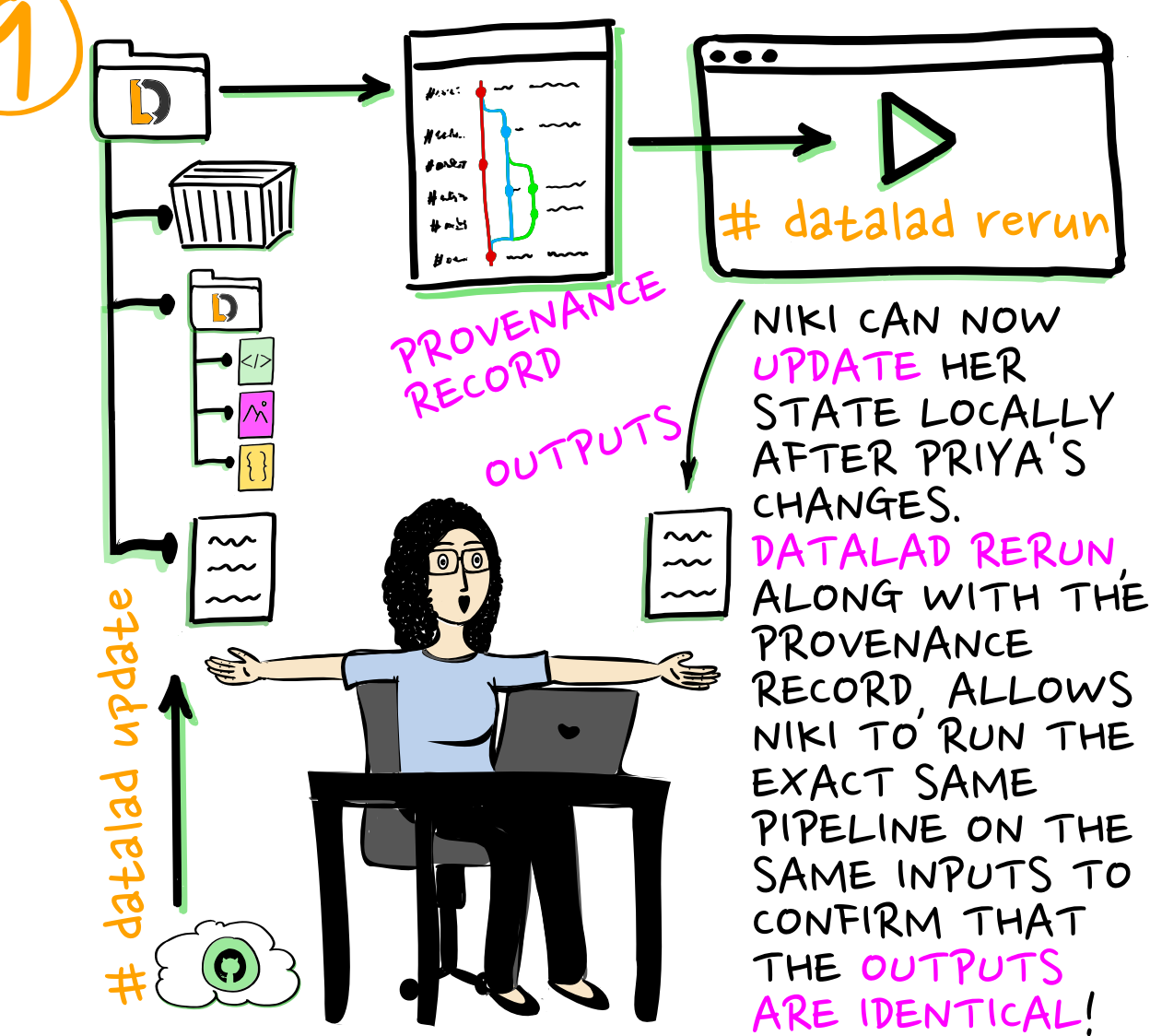


8

PRIYA WANTS TO RUN A REPRODUCIBLE PIPELINE ON THE DATA, TO SHARE WITH NIKI. SHE CONTAINERIZES THE PIPELINE, AND NESTS IT TOGETHER WITH THE INPUT DATASET INTO A NEW PROJECT DATASET:



9



NIKI CAN NOW UPDATE HER STATE LOCALLY AFTER PRIYA'S CHANGES. DATALAD RERUN ALONG WITH THE PROVENANCE RECORD, ALLOWS NIKI TO RUN THE EXACT SAME PIPELINE ON THE SAME INPUTS TO CONFIRM THAT THE OUTPUTS ARE IDENTICAL!

10



NIKI IS VERY SATISFIED WITH HER TEAM'S NEW DATA MANAGEMENT TOOL AND PRACTICES.

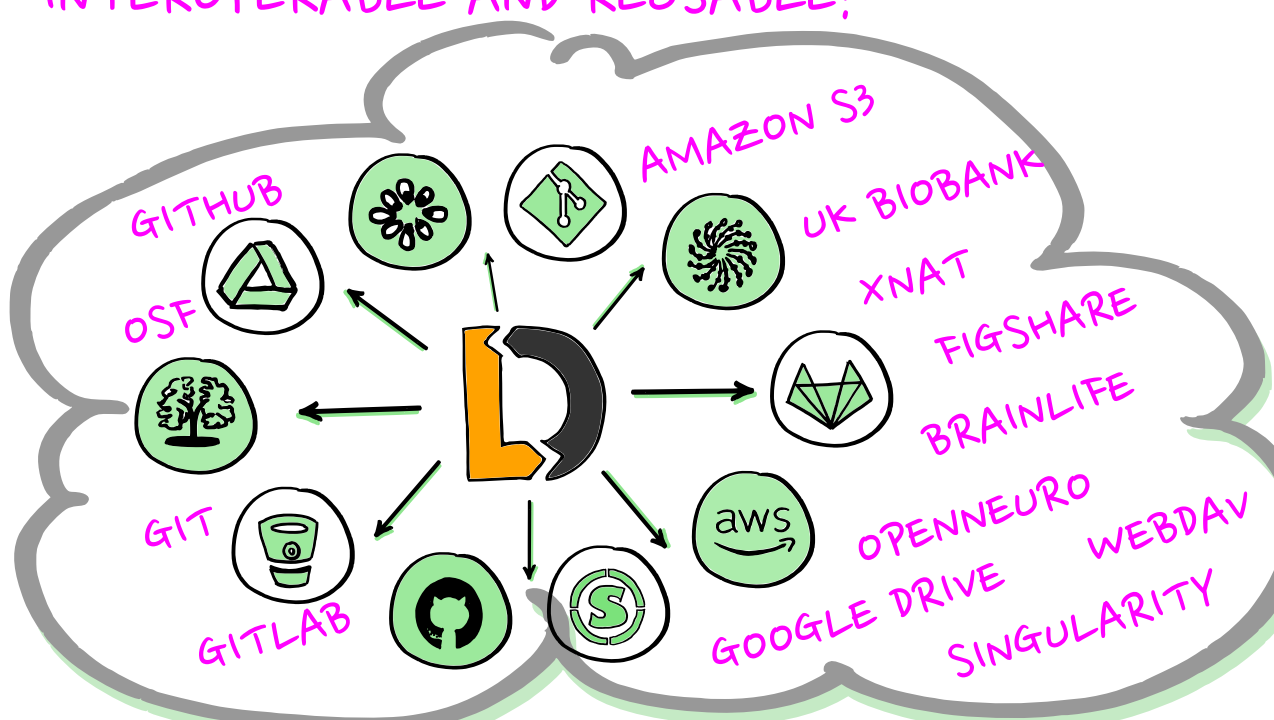
AS THEY EXPLORE MORE CAPABILITIES, SHE FINDS A WEALTH OF USAGE INFORMATION AND PRACTICAL EXAMPLES IN THE DATALAD HANDBOOK



<http://handbook.datalad.org/>

11

THEY ALSO LEARN THAT DATALAD'S INTEGRATIONS AND EXTENSIONS INCLUDE WIDELY USED DATA MANAGEMENT AND STORAGE SERVICES. THEY PLAN TO USE THESE WITH DATALAD TO MAKE THEIR DATA AND CODE MORE FINDABLE, ACCESSIBLE, INTEROPERABLE AND REUSABLE!



12

DATALAD RESOURCES

- datalad.org
- github.com/datalad
- info@datalad.org
- youtube.com/datalad
- twitter.com/datalad

JSH