

Machine Learning - Pipelines

- Open-source Version Control System for ML (DVC)
- Curious Containers (CC)

Jonas Annuscheit

09.04.2019
CBMI-Journal-Club

Inhalt

- Vokabular
- Problemstellung
- Data Version Control (DVC)
- Curious Containers (CC)
- DVC + CC
- Zusammenfassung

Inhalt

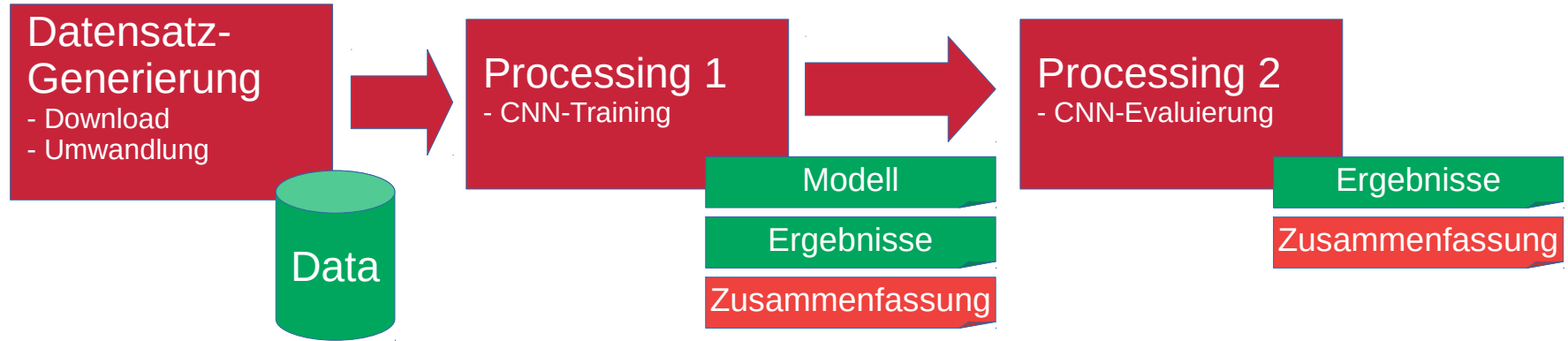
- Vokabular
- Problemstellung
- Data Version Control (DVC)
- Curious Containers (CC)
- DVC + CC
- Zusammenfassung

Vokabular

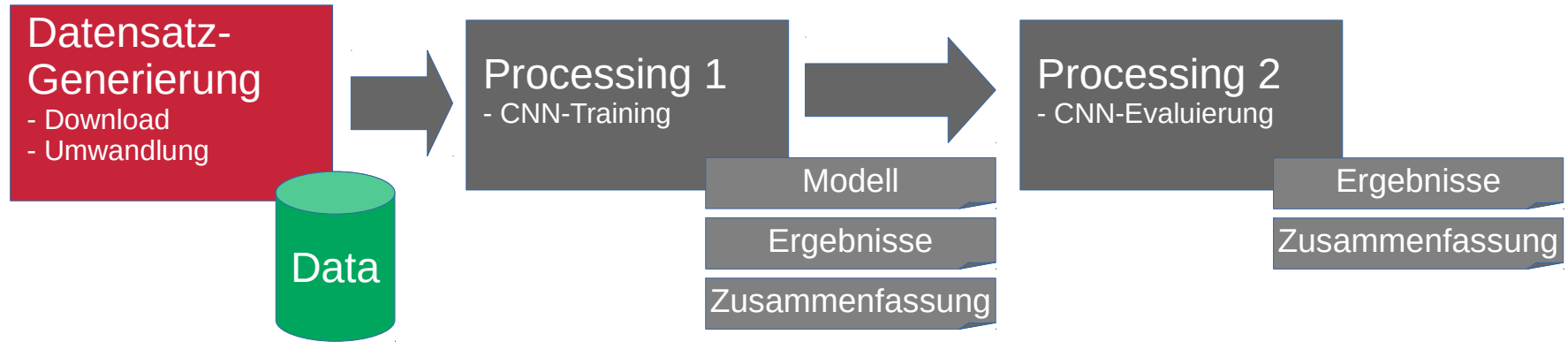
- **Projekt:** z.b.: Finden der besten Acc. für den CIFAR10-Datensatz.
- **Experiment:** An einer Stelle schrauben (z.B. Verschiedene Aktivierungsfunktion)
- **Sub-Experiment:** z.B.: Testen von 3, 5, 7 und 9 Schichten
- **Ungültiges Experiment:** Beinhaltet Programmierfehler

Problemstellung

■ Große Dateien
■ Kleine Dateien



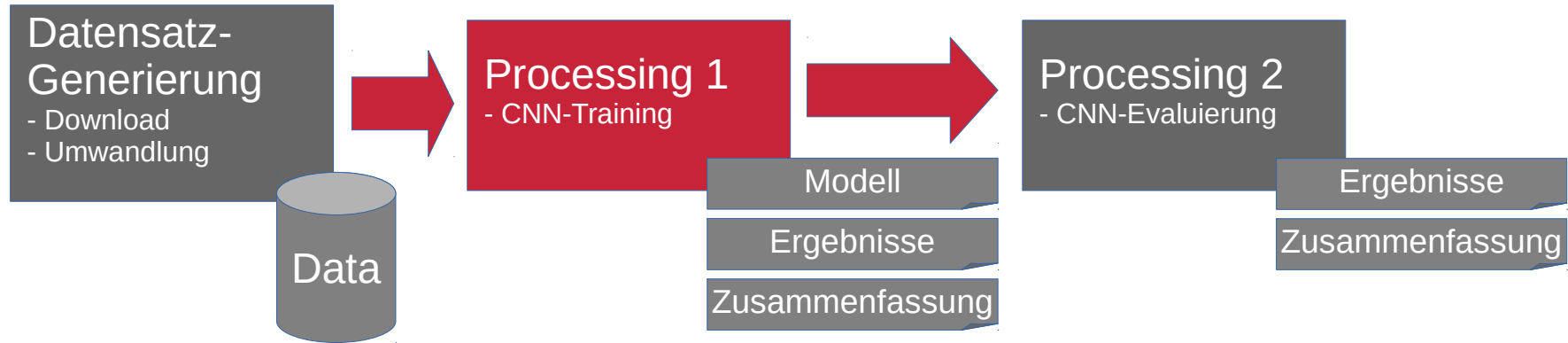
Problemstellung



Probleme:

- Verwerfung der Informationen, die zum Datensatz geführt haben.

Problemstellung

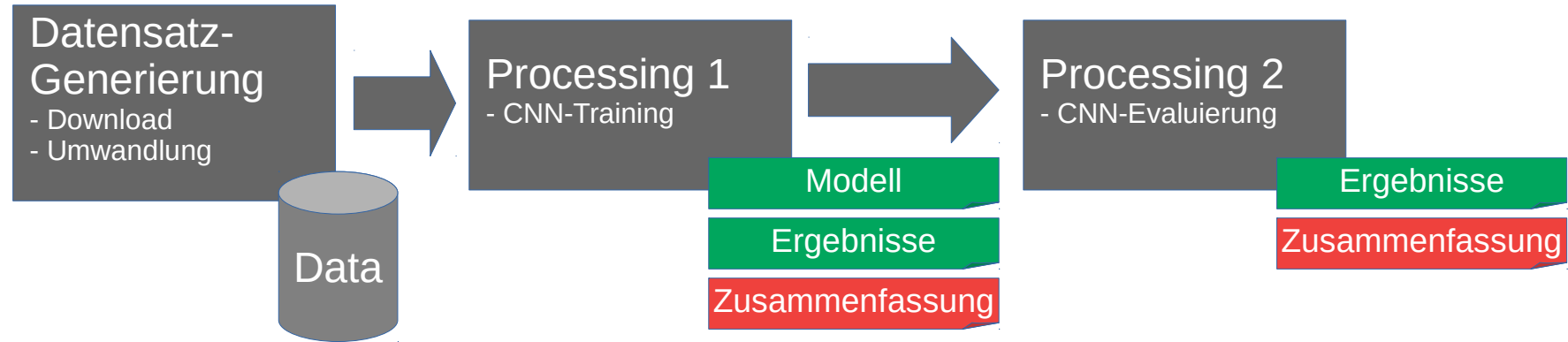


Probleme:

- Verwerfung der Informationen, die zum Datensatz geführt haben.
- Verwendete Parameter unbekannt
- Welche Datei zu welchem Ergebnis geführt hat, ist nicht nachvollziehbar.

Problemstellung

- Große Dateien
- Kleine Dateien



Probleme:

- Verwerfung der Informationen, die zum Datensatz geführt haben.
- Verwendete Parameter unbekannt
- Überschreibung der Dateien
- Welche Datei zu welchem Ergebnis geführt hat, ist nicht nachvollziehbar.
- Nichts sagende Benennung

Problemstellung

- **Reproduzierbarkeit:** Jeder kann die Ergebnisse Reproduzieren
- **Klarheit:** Leichte Verständlichkeit des Experiment-Aufbaues
- **Skalierbarkeit:** Mehrere Experimente parallel laufen lassen
- **Einfachheit:** Kein oder nur wenig extra Aufwand

Open-source Version Control System for ML (DVC)

Jonas Annuscheit

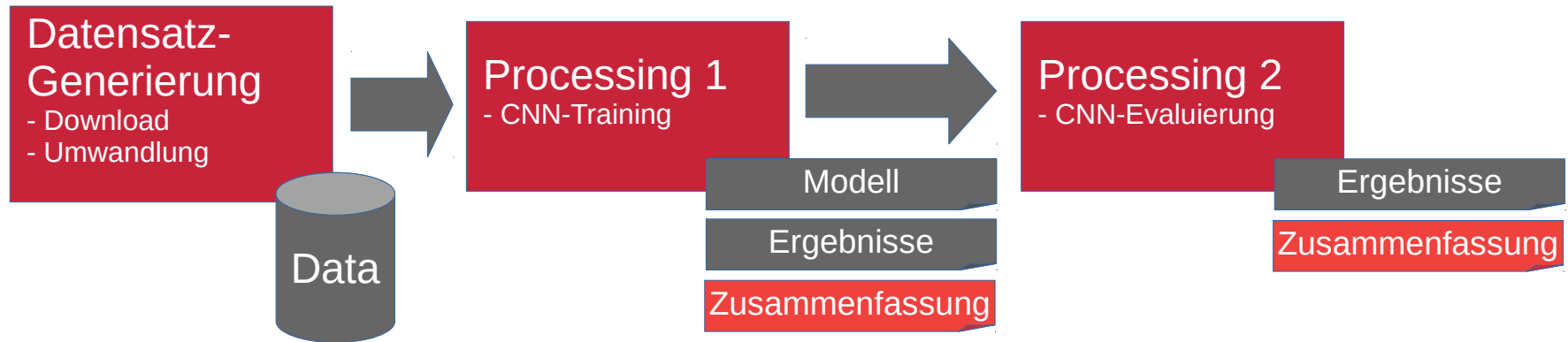
09.04.2019
CBMI-Journal-Club

DVC

https://www.youtube.com/watch?v=4h6l9_xeYA4

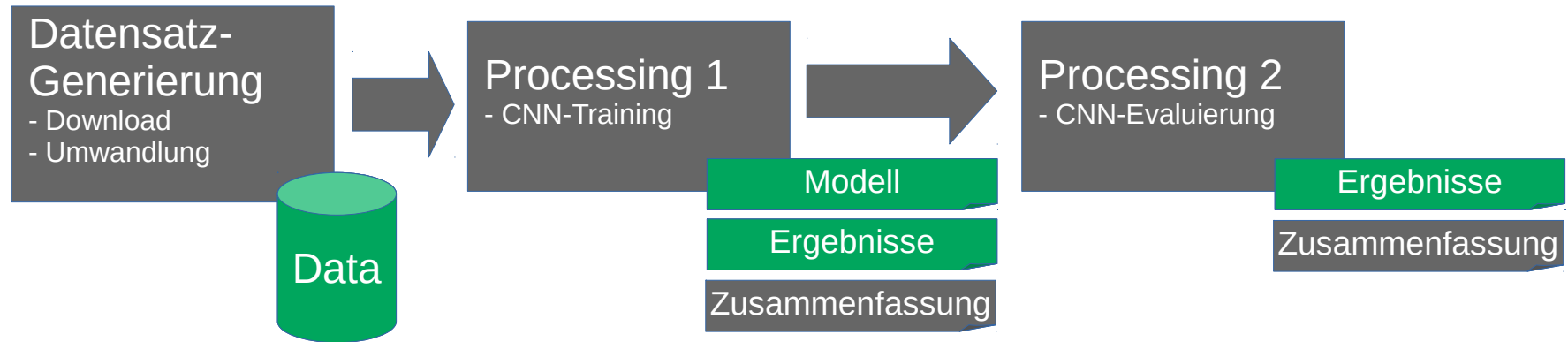
Erste Annahme

- Git ist eine gute Möglichkeit Source-Code zu verwalten



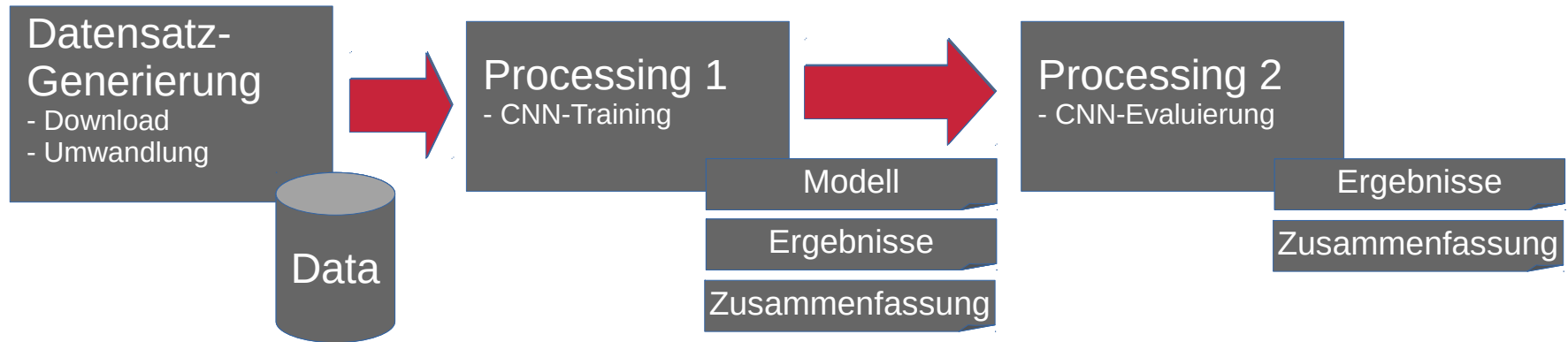
Zweite Annahme

- Große Dateien müssen extra gehandelt werden



Dritte Annahme

- Die Abfolge der Abrufe muss gespeichert werden.



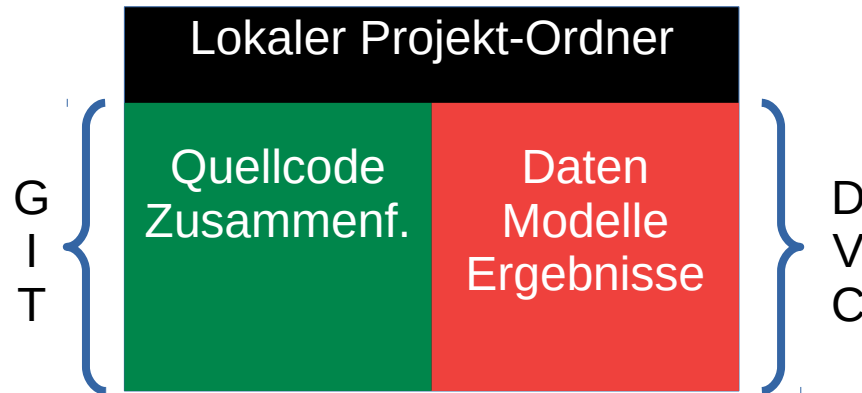
Zusammenfassung

Lokaler Projekt-Ordner

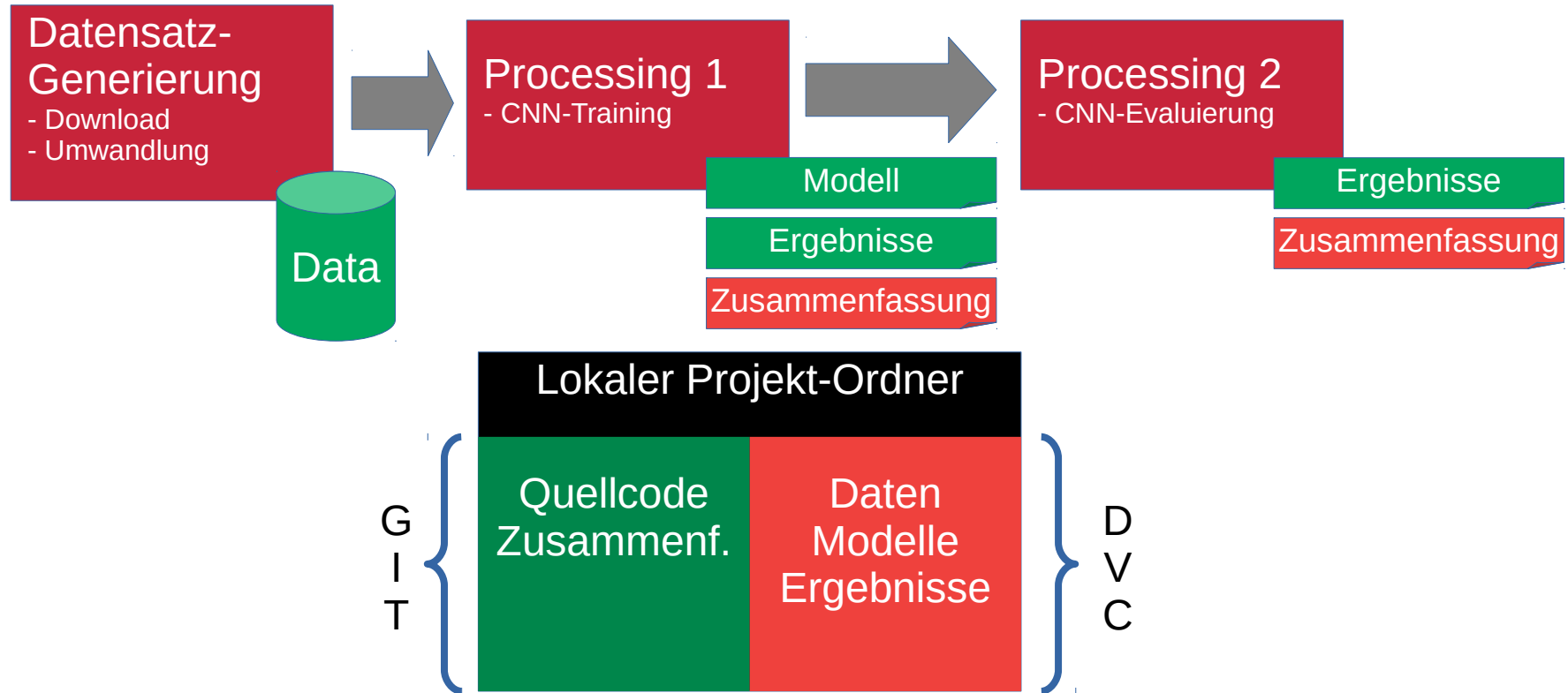
Quellcode
Zusammenf.

Daten
Modelle
Ergebnisse

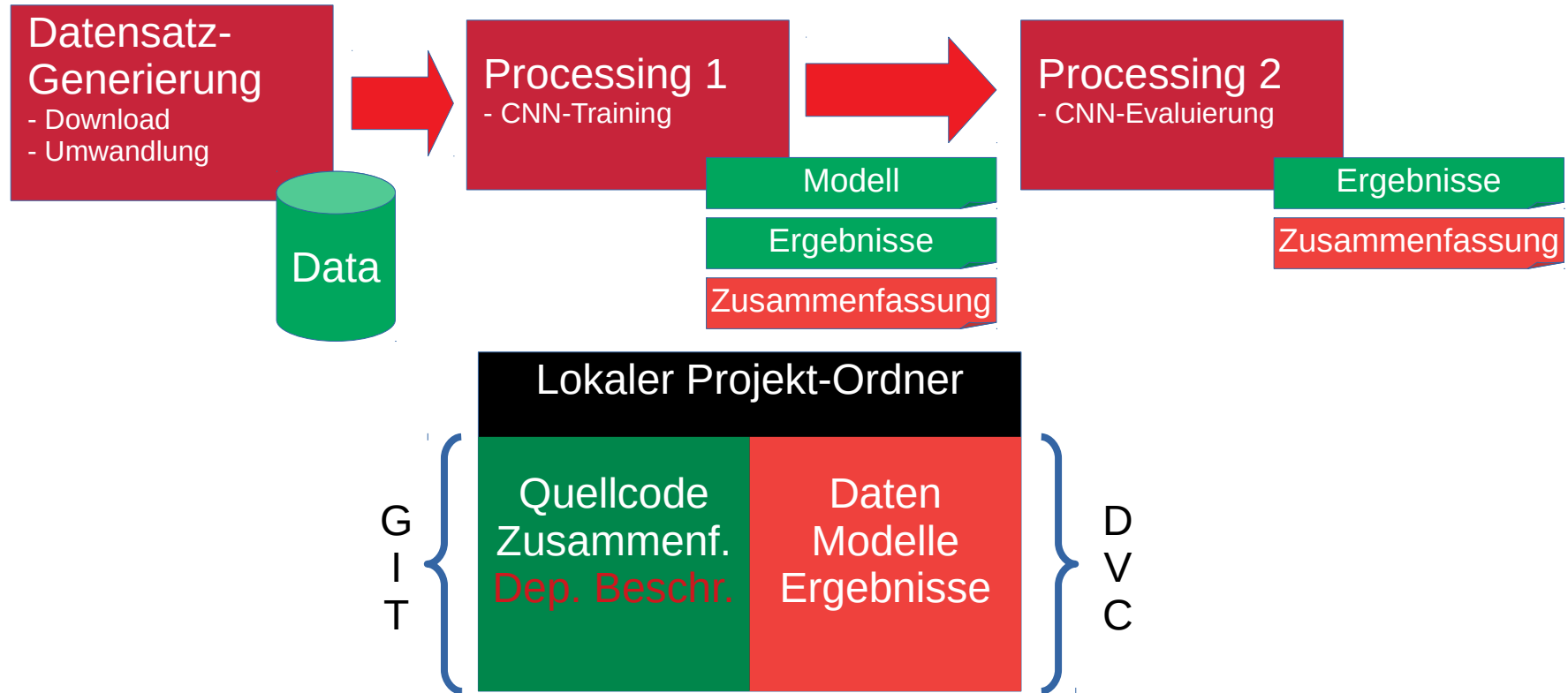
Zusammenfassung



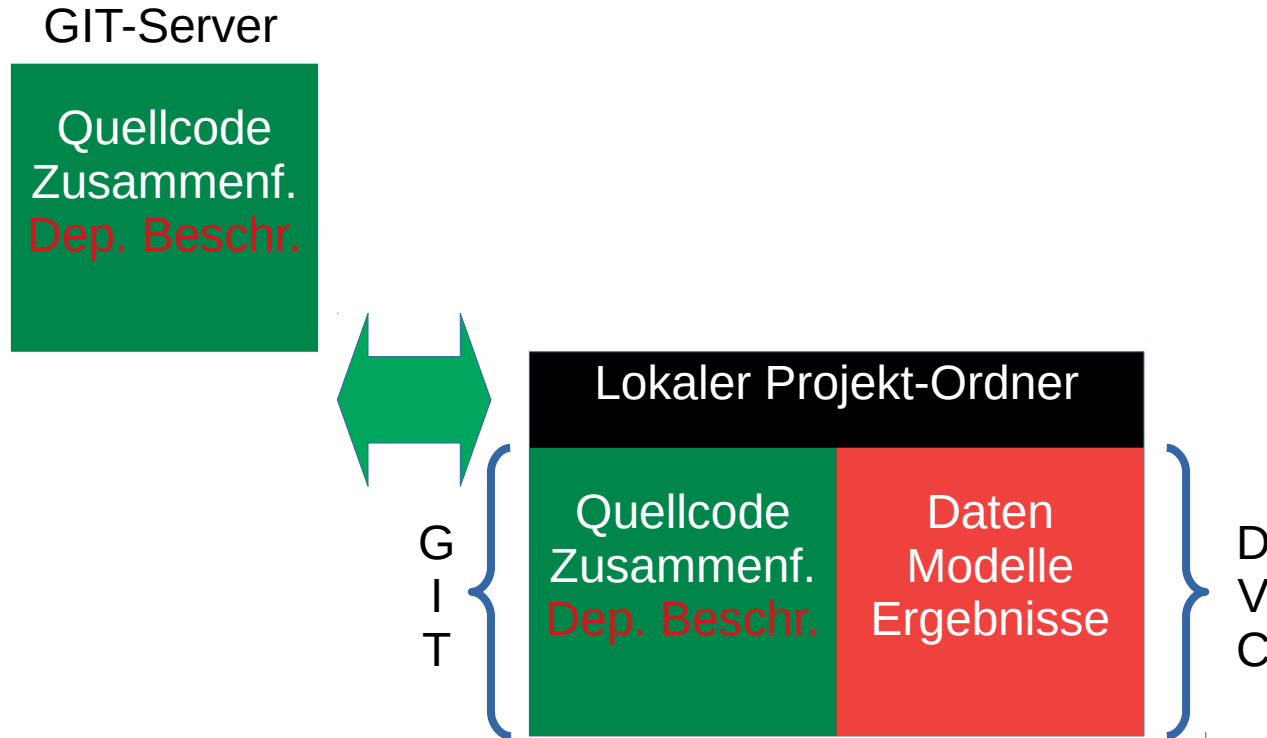
Zusammenfassung



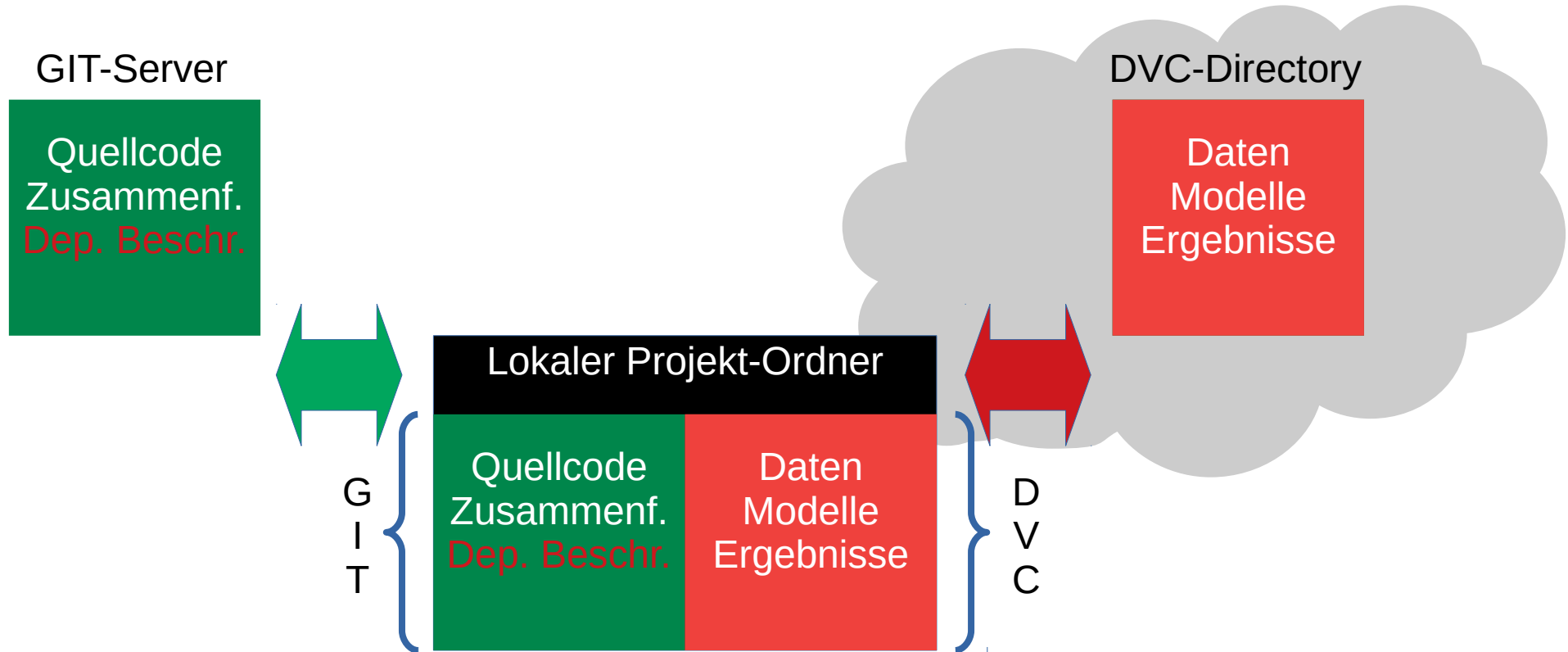
Zusammenfassung



Zusammenfassung



Zusammenfassung



DEMO: DVC

Komandos

```
1 # First you should install Anaconda and create an environment
2 conda create --name the_name_of_your_env pip git numpy appdirs decorator -y
3
4 # activate the environment
5 conda activate the_name_of_your_env    # in older versions: source activate the_name_of_your_env
6
7 # download source code
8 git clone https://github.com/mastaer/dvc-cc.git
9
10 # install the software that you need to run the code
11 pip install dvc-cc/dep/cc_core-7.0.0-py3-none-any.whl --ignore-installed    # soon this will be: pip install cc-core
12 pip install dvc-cc/dep/cc_faice-7.0.0-py3-none-any.whl --ignore-installed    # soon this will be: pip install cc-faice
13 pip install dvc-cc/dist/dvc_cc-0.1.0-py3-none-any.whl --ignore-installed
14
15 # GO to github or gitlab and create a new project: i.e. the project name: dvc_demo
16
17 # Clone the git: here is it the following link
18 git clone https://git.tools.f4.htw-berlin.de/annusch/dvc_demo.git
19
20 # Go the the project folder
21 cd dvc_demo
```

```
24 # CONVERT TO A DVC-Project,
25 # dvc init # This will be called by dvc-cc project create
26
27 # Create simple project
28 dvc-cc project create --mini_project
29
30 # instead of just calling "python create_some_data.py", we define with
31 # we use the --no-exec command, so that we can demonstrate later the dvc repro -P
32 dvc run -d create_some_data.py \
33         -o mydata.npy \
34         --no-exec \
35         -f _create_data.dvc \
36         python create_some_data.py
37
38 # we do this also for the train.py and for the eval.py
39 dvc run -d mydata.npy \
40         -d train.py \
41         -o train_model.npy \
42         -m train_metric.json \
43         --no-exec \
44         -f _train.dvc \
45         python train.py
46
47 dvc run -d mydata.npy \
48         -d train_model.npy \
49         -d eval.py \
50         -m test_metric.json \
51         -f eval.dvc \
52         --no-exec
53         python eval.py
```

```
56 # NOW we can run the experiment
57 dvc repro -P
58
59 # If we delete one file, for example the test_metric.json and call dvc repro -P again, it will only run the stage that
   is needed
60 rm test_metric.json
61 dvc repro -P
62
63
64 # show the dependencies
65 dvc pipeline show --ascii eval.dvc
66
67 # show a metric that
68 dvc metrics show -t json -x L1_Test -a
69
70
71
72
73 # if you are finish with testing, and this is nothing for you, you can remove the env with:
74 conda deactivate
75 conda remove -n the_name_of_your_env --all -y
__
```


Problemstellung (DVC)

- **Reproduzierbarkeit:** (fast) jeder kann die Ergebnisse reproduzieren
- **Klarheit:** Leichte Verständlichkeit des Experiment-Aufbaues
- **Skalierbarkeit:** Mehrere Experimente parallel laufen lassen
- **Einfachheit:** Kein oder nur wenig extra Aufwand

Curious Containers (CC)

Jonas Annuscheit

09.04.2019
CBMI-Journal-Club

Informationen

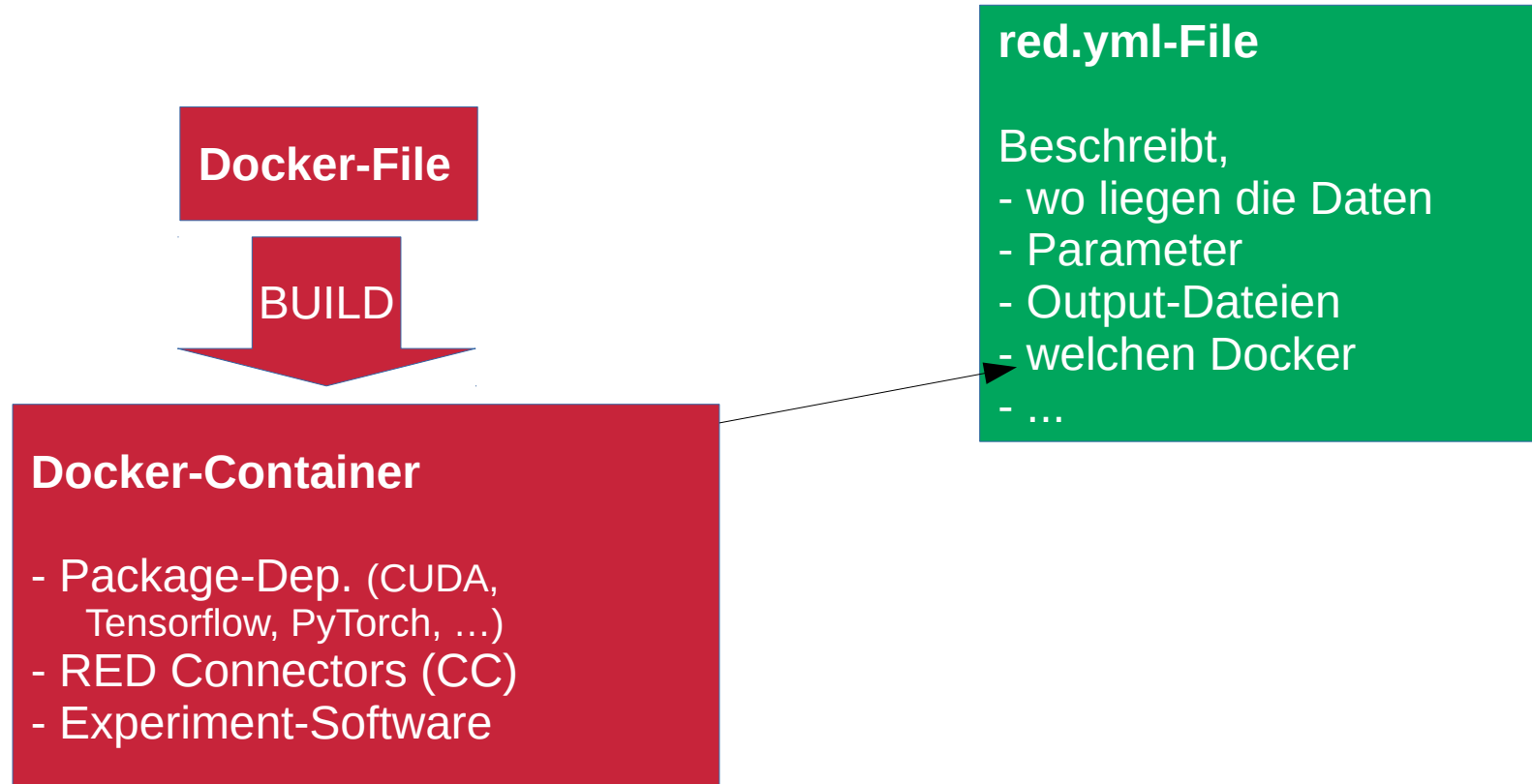
- **Autoren:** Christoph Jansen & Bruno Schilling
- **Ziele:** Reproduzierbarkeit & Skalierbarkeit
- **Paper:** Christoph Jansen et al. Reproducibility and Performance of Deep Learning Experiments for Cancer Detection in Pathological Images

Informationen

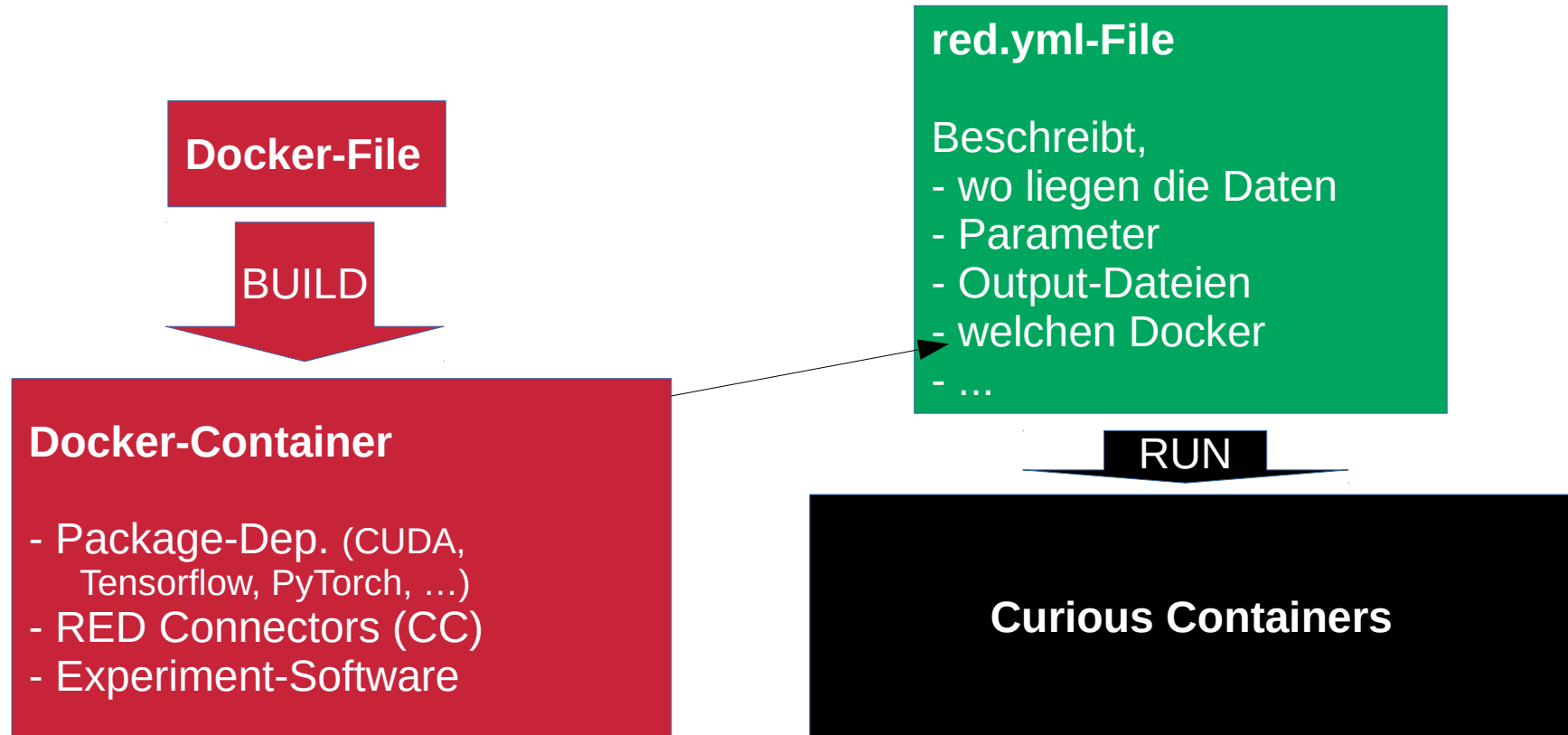
- **Autoren:** Christoph Jansen & Bruno Schilling
- **Ziele:** Reproduzierbarkeit & Skalierbarkeit
- **Paper:** Christoph Jansen et al. Reproducibility and Performance of Deep Learning Experiments for Cancer Detection in Pathological Images



Die Idee: Ein Experiment = Ein Docker-Container + red.yml-File



Die Idee: Ein Experiment = Ein Docker-Container + red.yml-File



Problemstellung (CC)

- **Reproduzierbarkeit:** Jeder kann die Ergebnisse reproduzieren
- **Klarheit:** Leichte Verständlichkeit des Experiment-Aufbaues
- **Skalierbarkeit:** Mehrere Experimente parallel laufen lassen
- **Einfachheit:** Kein oder nur wenig extra Aufwand

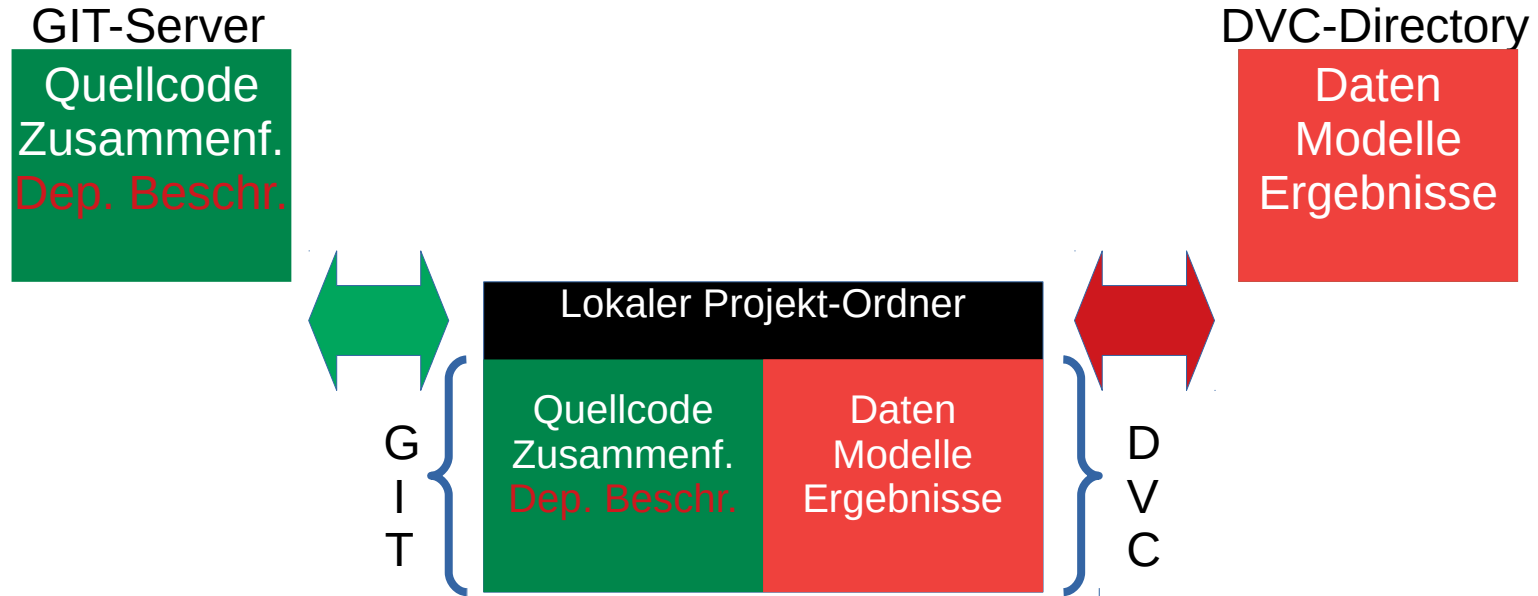
DVC + CC

Jonas Annuscheit

09.04.2019
CBMI-Journal-Club

Idee: Schritt 1

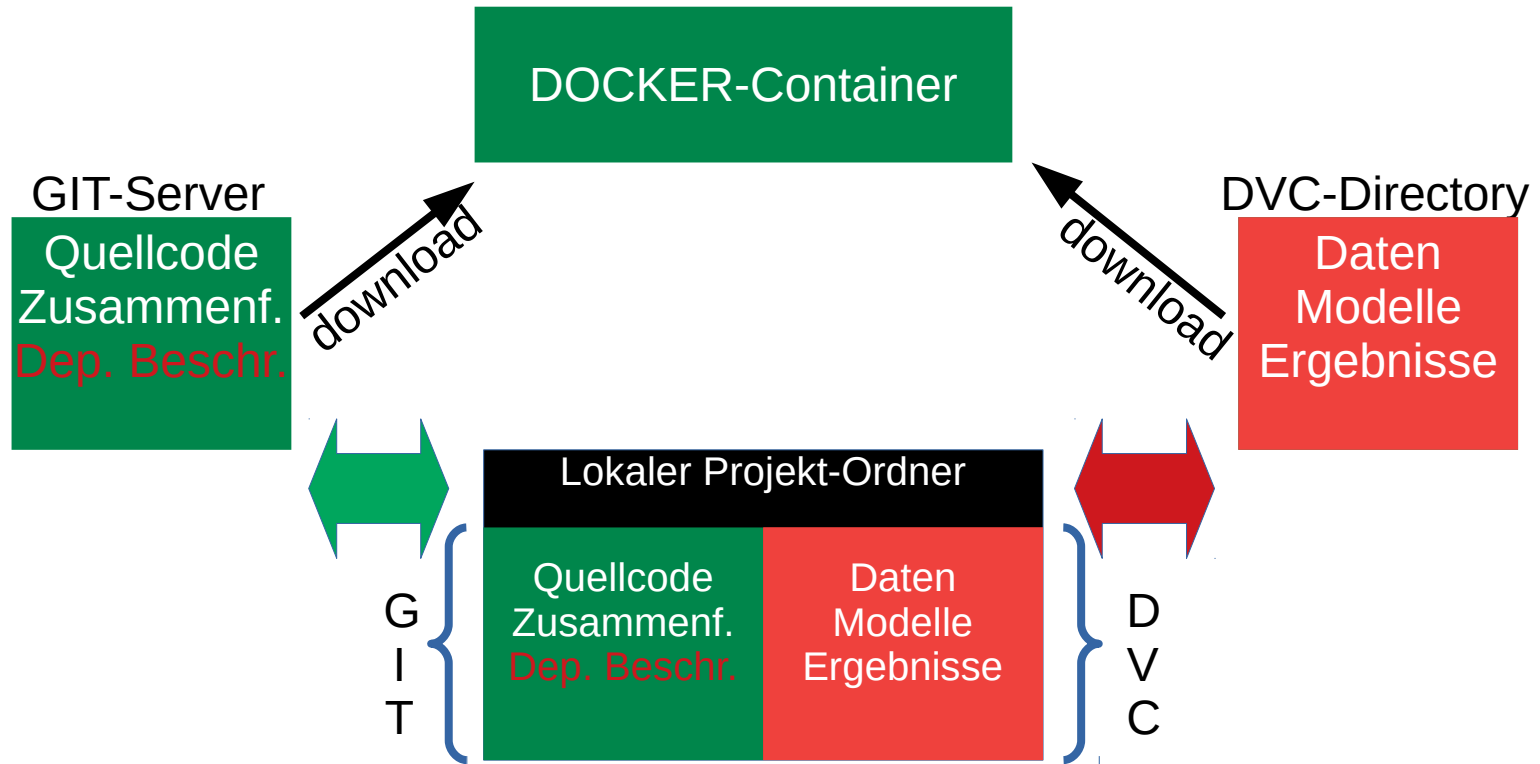
- Verwendung von DVC



Idee: Schritt 1.5

- Selber Docker-Container für alle Projekte
- Docker beinhaltet statt Source-Code vom Experiment, GIT und DVC-Repro.-Downloader

Idee: Schritt 1.5



Idee: Schritt 2

- red.yml-File aus den Informationen des GIT bzw. DVC automatisiert generieren

DEMO: DVC

Komandos

```
1 # First you should install Anaconda and create an environment
2 conda create --name the_name_of_your_env pip git numpy appdirs decorator -y
3
4 # activate the environment
5 conda activate the_name_of_your_env # in older versions: source activate the_name_of_your_env
6
7 # download source code
8 git clone https://github.com/mastaer/dvc-cc.git
9
10 # install the software that you need to run the code
11 pip install dvc-cc/dep/cc_core-7.0.0-py3-none-any.whl --ignore-installed # soon this will be: pip install cc-core
12 pip install dvc-cc/dep/cc_faice-7.0.0-py3-none-any.whl --ignore-installed # soon this will be: pip install cc-faice
13 pip install dvc-cc/dist/dvc_cc-0.1.0-py3-none-any.whl --ignore-installed
14
15 # GO to github or gitlab and create a new project: i.e. the project name: dvc_demo
16
17 # Create a Folder on a server for the extern dvc directory
18
19 # Clone the git:
20 git clone https://git.tools.f4.htw-berlin.de/annusch/demo_dvc.git|
21
22 # Go the the project folder
23 cd demo_dvc
```

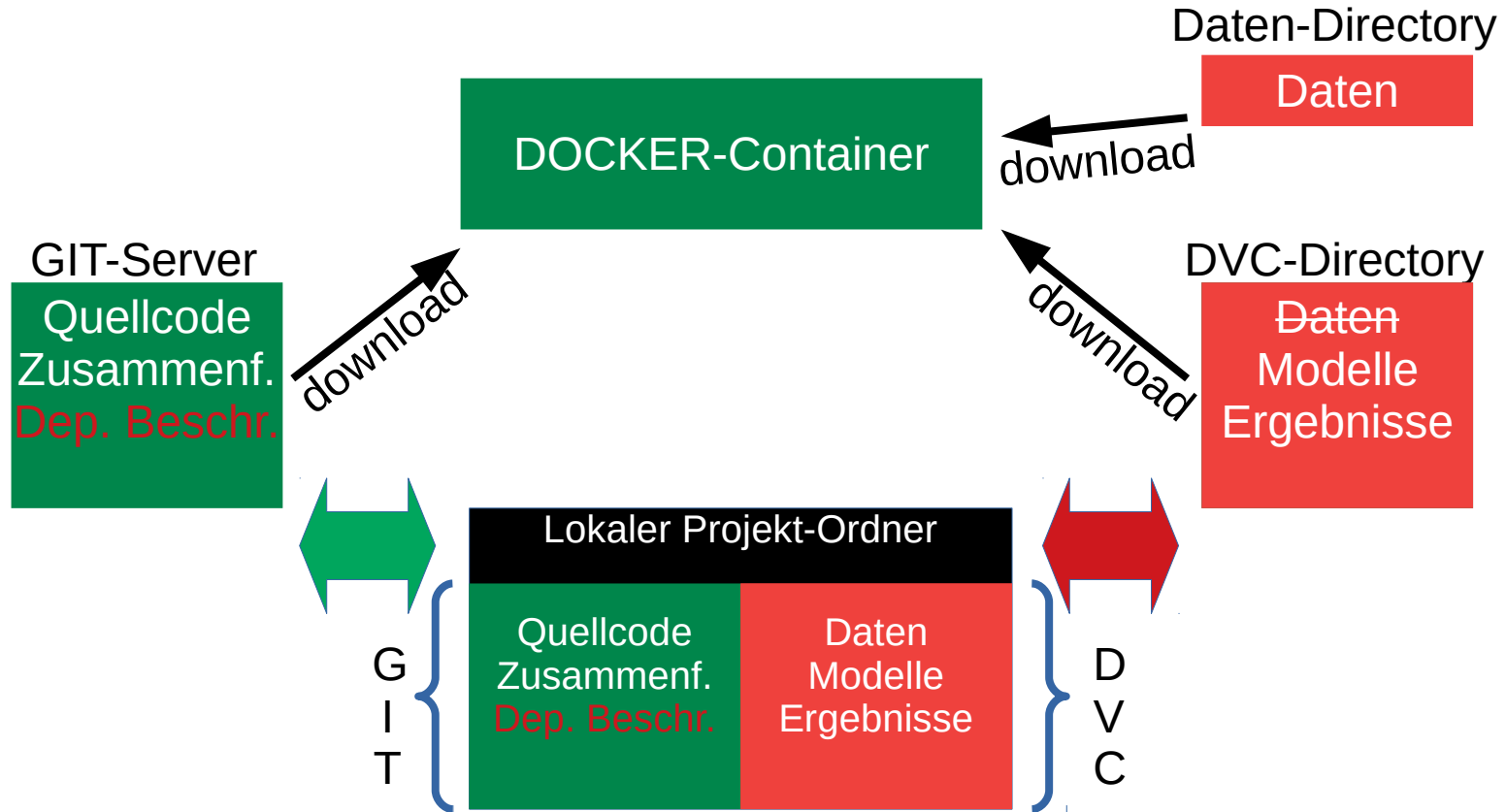
```
26 # CONVERT TO A DVC-Project,
27 # dvc init # This will be called by dvc-cc project create
28
29 # Create simple project
30 dvc-cc project create
31
32 # create a .dvc/config.local file: It should look like this (nano .dvc/config.local)
33 ['remote "nas"']
34 url = ssh://annusch@avocado01.f4.htw-berlin.de/data/ldap/jonas/demo_dvc
35 ask_password = true
36
37 [core]
38 remote = nas
39
40 # create the dataset
41 dvc repro _generate_data.dvc
42
43 # push everything
44 dvc-cc git commit_and_push
```

```
46 # build a experiment in a branch
47 dvc-cc git branch first_experiment
48
49 # add the following line to nano code/argparser.py
50     parser.add_argument('-n', '--num_of_samples', default=10)
51 # edit line 26 from code/train.py to
52 num_of_samples = int(args.num_of_samples)
53
54 # Create some sub-experiments:
55 dvc-cc project dummy -p "-n 2"
56 dvc-cc project dummy -p "-n 5"
57 dvc-cc project dummy -p "-n 10"
58 dvc-cc project dummy -p "-n 20"
59 dvc-cc project dummy -p "-n 100"
60
61 dvc-cc git commit_and_push
62
63 # create red yml file
64 dvc-cc red add_job -T
65
66 # run the job
67 dvc-cc jobs run
68
```


Problemstellung (DVC + CC)

- **Reproduzierbarkeit:** (fast-fast-fast) jeder kann die Ergebnisse reproduzieren
- **Klarheit:** Leichte Verständlichkeit des Experiment-Aufbaues
- **Skalierbarkeit:** Mehrere Experimente parallel laufen lassen
- **Einfachheit:** Kein oder nur wenig extra Aufwand

What is next?



Zusammenfassung (Vorschläge)

Jonas Annuscheit

09.04.2019
CBMI-Journal-Club

Vorschläge

- mehrere Experimenten => DVC
- Ein Experiment == Ein Branch
- Automatisierte Experiment-Namen-Vergabe
(z.B. Branch-Name + Parameters) & Benennung aller Dateien danach
- Verwendung von Clustern => Nutzung von CC
mit DVC_CC

DANKE

Jonas Annuscheit

09.04.2019
CBMI-Journal-Club