

Poisoned Pickles Make You ill

Adrian Gonzalez-Martin

Head of ML Serving

agm@seldon.io



About me



Adrian Gonzalez-Martin
Head of ML Serving, Seldon
Fellow, The Institute for Ethical AI & ML

@adriangonz@mastodon.social

github.com/adriangonz

Security Challenges in Machine Learning

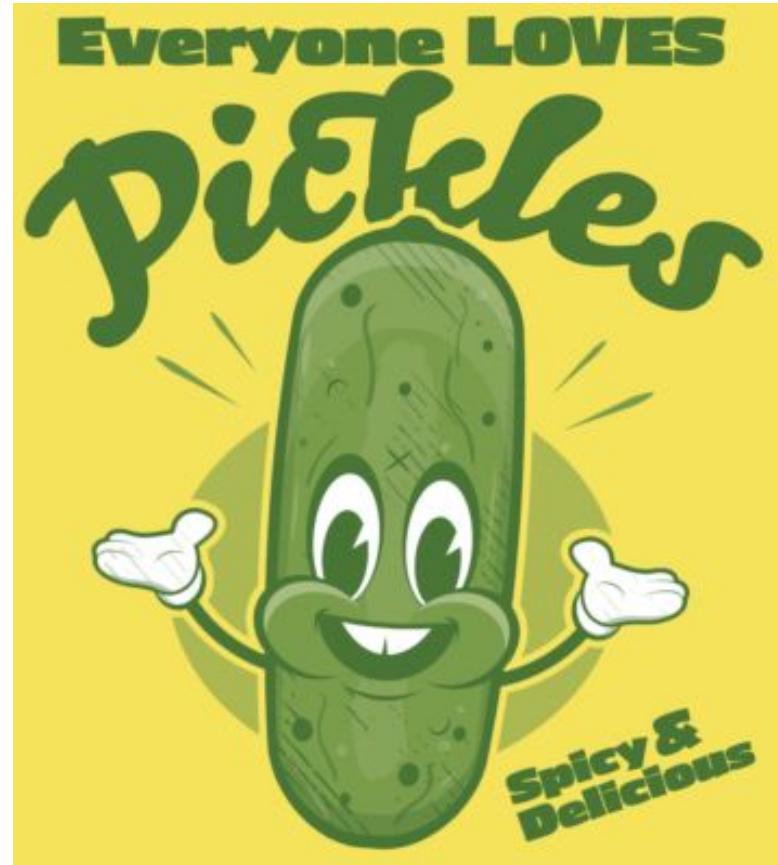
- Security challenges in devops and software space more known
- Key challenges in machine learning still being defined as best practice
- Impossible to make systems **unhackable**, but possible to **mitigate undesired outcomes**

The solutions will be technical in nature but will ultimately still rely on humans and process.

Everyone loves Pickles!

Everyone loves pickles!

- Pickles are Python's native **serialisation format**
- **Any Python object** can be serialised into a Pickle
 - ◆ Including functions and **arbitrary Python code**
- Hugely popular in the ML space



Pickles are everywhere

→ Major ML libraries rely on Pickles (i.e. they are exposed to Pickle's issues)

The screenshot displays three separate sections illustrating the use of pickles in machine learning frameworks:

- Scikit-learn (Left):** A sidebar for "Model persistence" under version 1.0.2. It includes links for "Prev", "Up", "Next", and "Other versions". A note says "Please cite us if you use the software." Below, it lists "9. Model persistence" with "9.1. Python specific serialization" and "9.2. Interoperable formats".
- PyTorch (Top Right):** The PyTorch documentation page for `torch.load`. The header includes "Get Started", "Ecosystem", "Mobile", "Blog", "Tutorials", "Docs", and "Resources". The main content shows code snippets for `dump` and `load`, notes about file-like objects, and a section on security/maintainability limitations.
- TensorFlow (Bottom Right):** The TensorFlow API documentation for `TORCH LOAD`. It shows the function signature `torch.load(f, map_location=None, pickle_module=pickle, **pickle_load_args) [SOURCE]` and a note explaining its behavior, mentioning Python's unpickling facilities and storage mappings.

Below the main sections, there is a large block of Python code demonstrating the serialization and deserialization of a custom layer between TensorFlow and PyTorch:

```
layer = CustomLayer(5)
layer.var.assign(2)

serialized_layer = keras.layers.serialize(layer)
new_layer = keras.layers.deserialize(
    serialized_layer, custom_objects={"CustomLayer": CustomLayer}
)
```

Security Risks with Pickles

By design, Pickles can serialise **any Python object**

Security Risks with Pickles II

Including **arbitrary code executions!**

```
[50]: import joblib
       model_safe = joblib.load("fml-artifacts/safe/model.joblib")

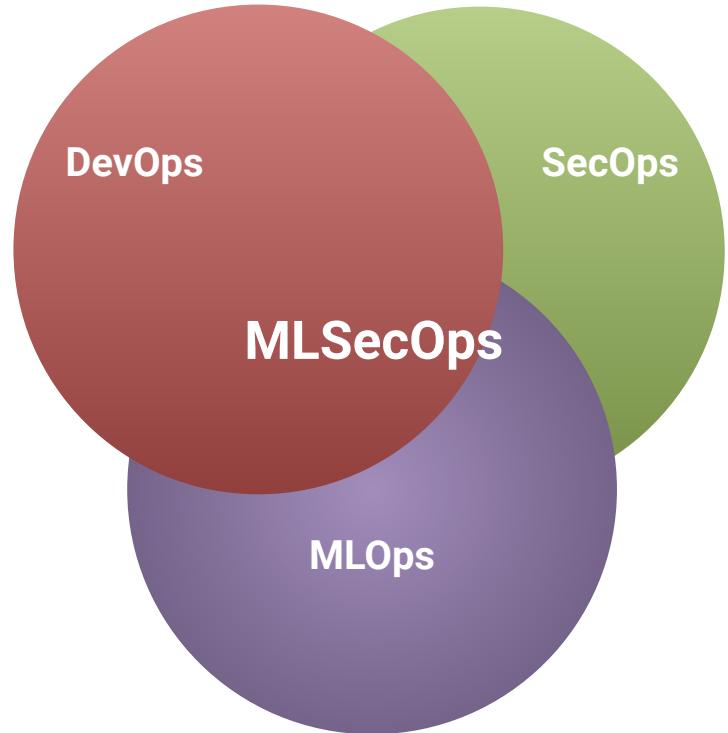
[51]: model_safe.predict(X[:1])
[51]: array([0])

[1]: import types, os, base64
      def __reduce__(self):
          # This is basically base64 for cmd = "env > pwnd.txt"
          cmd = base64.b64decode("ZW52ID4gcHduZC50eHQ=").decode()
          return os.system, (cmd,)

[2]: model_safe.__class__.__reduce__ = types.MethodType(__reduce__, model_safe.__class__)
      ...
[ ]: !mkdir -p fml-artifacts/unsafe/
[ ]: joblib.dump(model_safe, "fml-artifacts/unsafe/model.joblib")
[3]: with open("fml-artifacts/unsafe/model.joblib", "rb") as f: print(f.readlines())
[3]: [b'\x80\x04\x95]\x00\x00\x00\x00\x00\x00\x8c\x05posix\x94\x8c\x06system\x94\x93\x94\x8c\x0eenv > pwnd.txt\x94\x85\x94R\x94.]
```

What is MLSecOps?

- The Extension of DevOps and Security with machine learning infrastructure as first class citizen
- The Intersection of ML infrastructure, developer operations / automation, and security policies



Trusted AI Principles

The 8 LFAI Principles for Trusted AI – (R)REPEATS

Reproducibility

Robustness

Equitability

Privacy

Explainability

Accountability

Transparency

Security

DLF AI & DATA

02/10/2021

4

LFAI MLSecOps Top 10 for Machine Learning

[Flawed ML Security talk from KubeCon 2022](#)



#	MLSecOps Top 10
1	Unrestricted Model Endpoints
2	Access to Model Artifacts
3	Artifact Exploit Injection
4	Insecure ML Systems/Pipeline Design
5	Data & Infra Misconfigurations
6	Supply Chain Vulnerabilities in ML Code
7	IAM & RBAC Failures for ML Services
8	ML Infra / ETL / CI / CD Integrity Failures
9	Observability, Reproducibility & Lineage
10	ML-Server Side Request Forgery

Other issues with Pickles

Besides security issues, **Pickles also suffer from other issues**:

→ **Python version incompatibility** issues

- ◆ Pickles generated in one version may not work with other versions

→ **Framework version incompatibility** issues

- ◆ No way to version framework in Pickles

→ Hard to **troubleshoot** the above

e.g. error from using the wrong sklearn version



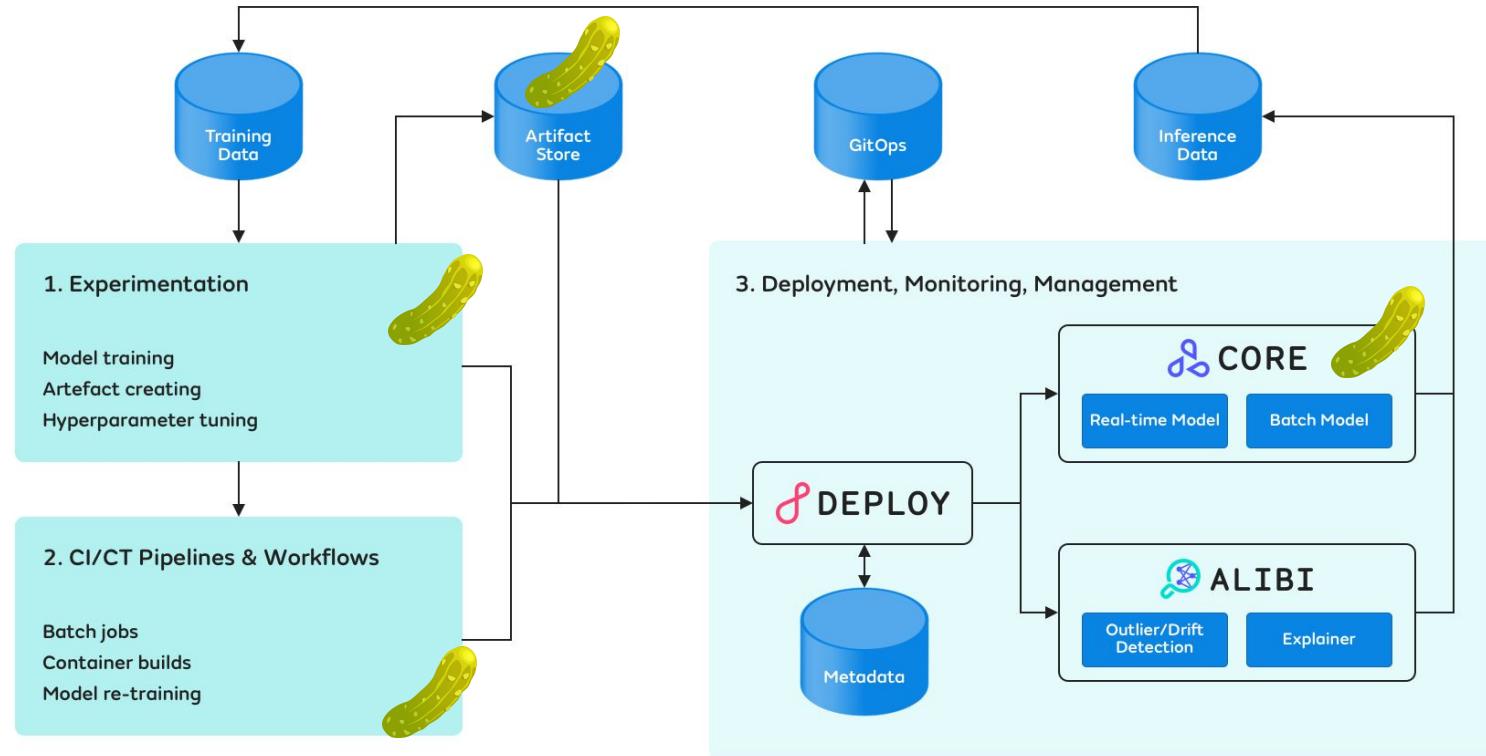
```
File "/opt/conda/lib/python3.8/site-packages/joblib/numpy_pickle.py", line 577, in _unpickle
    obj = unpickler.load()
File "/opt/conda/lib/python3.8/pickle.py", line 1212, in load
    dispatch[key[0]](self)
File "/opt/conda/lib/python3.8/site-packages/joblib/numpy_pickle.py", line 402, in load_build
    Unpickler.load_build(self)
File "/opt/conda/lib/python3.8/pickle.py", line 1705, in load_build
    setstate(state)
File "sklearn/tree/_tree.pyx", line 714, in sklearn.tree._Tree.__setstate__
File "sklearn/tree/_tree.pyx", line 1418, in sklearn.tree._Tree._check_node_ndarray
ValueError: node array from the pickle has an incompatible dtype:
- expected: {'names': ['left_child', 'right_child', 'feature', 'threshold', 'impurity', 'n_node_samples', 'weighted_n_node_samples', 'missing_value'], 'formats': ['<1b', '<1b', '<f8', '<f8', '<1b', '<f8', 'u1'], 'offsets': [0, 8, 16, 24, 32, 40, 48, 56], 'itemsize': 64}
- got : [(['left_child', '<1b'], ('right_child', '<1b'), ('feature', '<1b'), ('threshold', '<f8'), ('impurity', '<f8'), ('n_node_samples', '<1b'), ('weighted_n_node_samples', '<f8'))]
```

Poisoned Pickles make you ill

Some background knowledge...

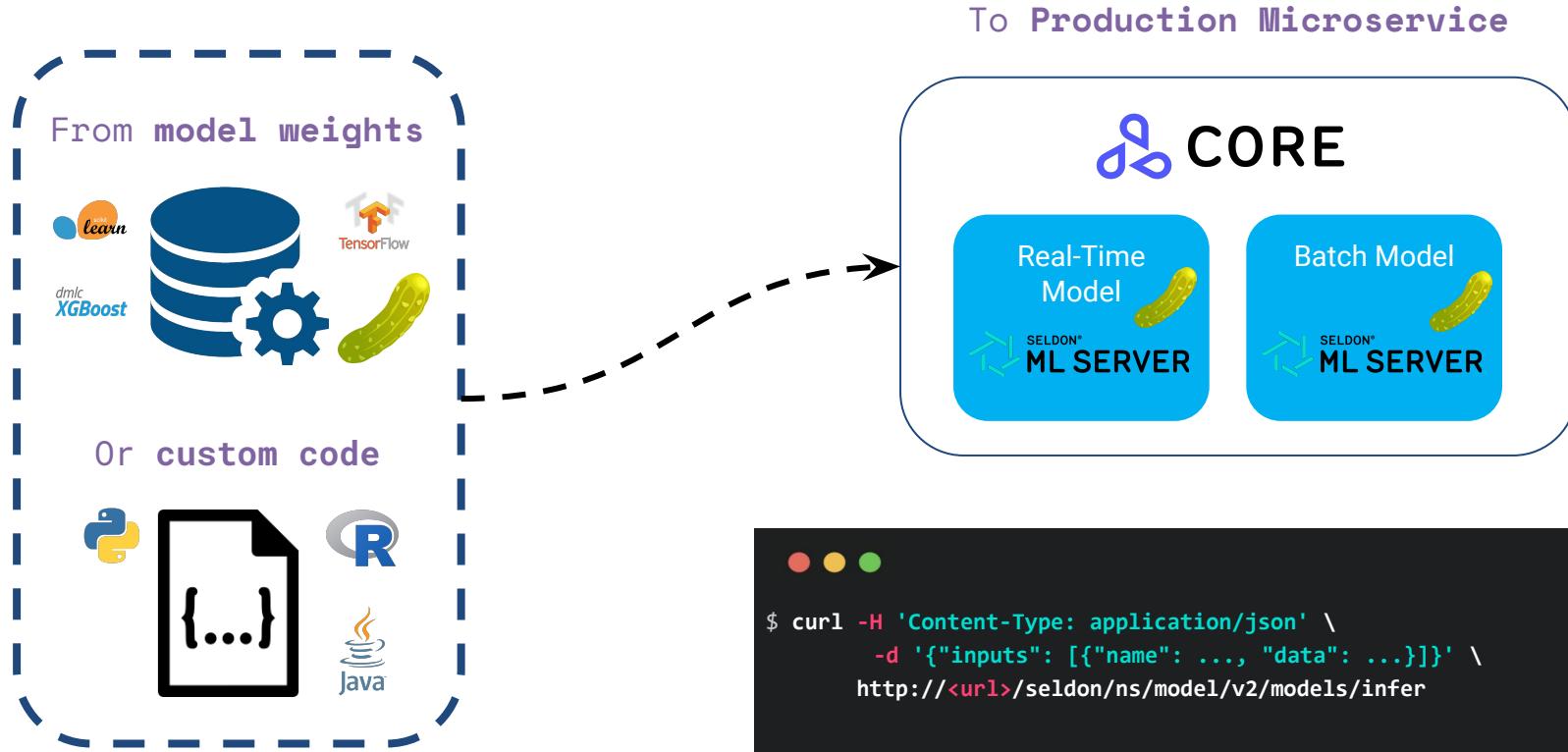


MLOps Serving Reference Architecture

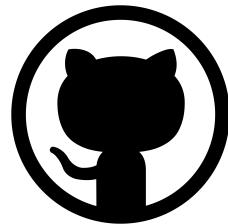


To learn more about production MLOps, catch up [Alejandro's talk on the State of Production ML](#)

TLDR; Machine Learning Deployed ✓



Demo



<https://github.com/adriangonz/sigstore-talk>

Poisoned Pickles Make You ill

- It's incredibly **easy to poison** a Pickle
- It's incredibly **hard to detect** if a Pickle has been poisoned

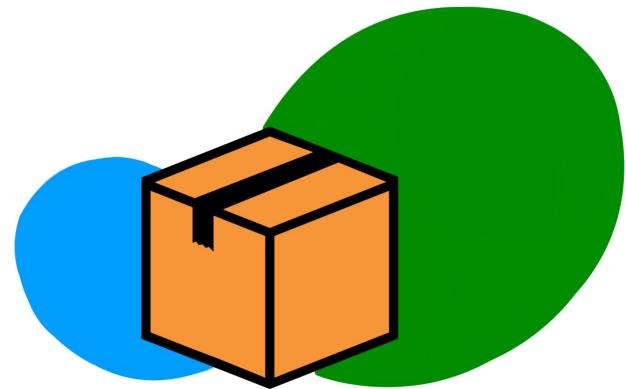


<https://huggingface.co/docs/hub/security-pickle>

Ok... Pickles may not be that great

What can we do? Can we have "*higher-quality Pickles*"?

- **Don't use Pickle** (e.g. use ONNX instead)
- Use tools like **SKOps** to mitigate Pickle risk
[see talk by Merve Noyan on PyData 2022]



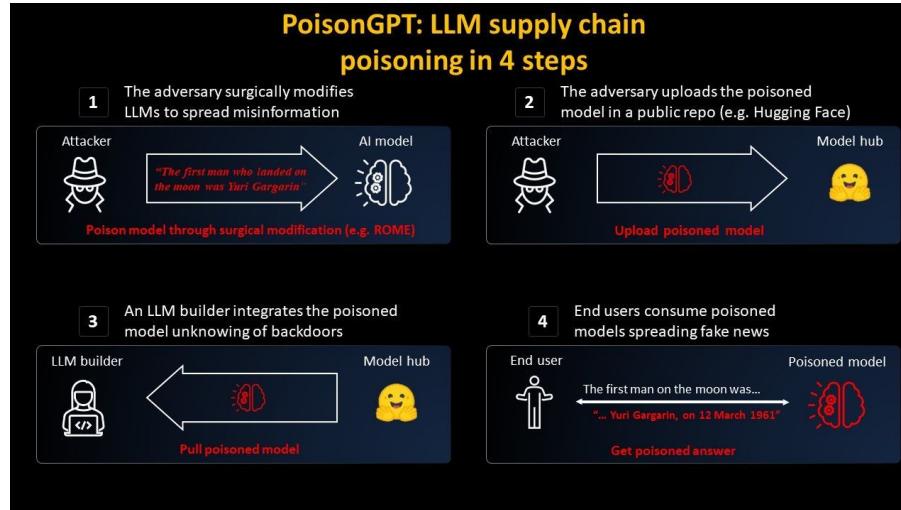
<https://skops.readthedocs.io/en/stable/>

Ok... Pickles may not be that great

Even with "higher-quality Pickles", there are risks....

PoisonGPT:

<https://blog.mithrilsecurity.io/poisongpt-how-we-hid-a-lobotomized-lm-on-hugging-face-to-spread-fake-news/>



Tamper-proof Pickle jars

Tamper-proof Pickle jars

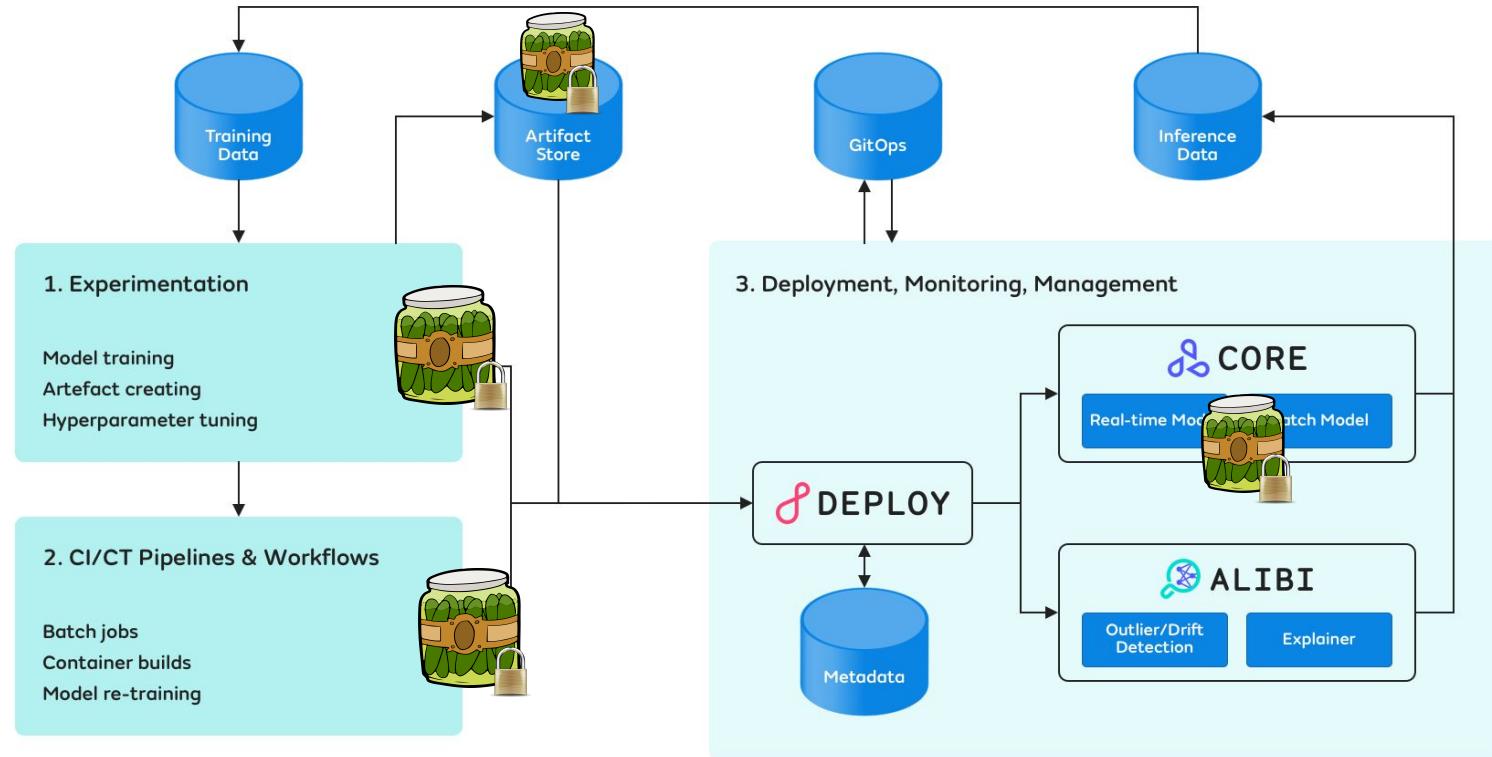
There's always a risk of an attacker tampering with the artefact:

- Security requirements are analogous to application development, which **cannot expect to be solved with just magical scanning**
- We need "**trust or discard**" mechanisms that need to be in play through the automation touchpoints of CI / CD / ETL

The screenshot shows the MITRE ATLAS website. The top navigation bar includes links for Home, Techniques, Matrices, Navigator, Tactics, Techniques, Case Studies, and Resources. The main content area has a sidebar titled "Techniques" with sections for ATLAS (Reconnaissance, Resource Development, Initial Access), ML Supply Chain Compromise, GPU Hardware, ML Software, Data, and Model. The main panel displays the "ML Supply Chain Compromise: Model" page, which includes a "Summary" section describing how machine learning systems often rely on open-sourced models and how they can be compromised. It also lists related IDs, tactics, parent techniques, case studies, and other subtechniques.

Catalogue of **attack vectors into ML systems**: <https://atlas.mitre.org/>

Tamper-proof Pickle jars



Tamper-proof Pickle jars

Can't we just *sign* the artefact?

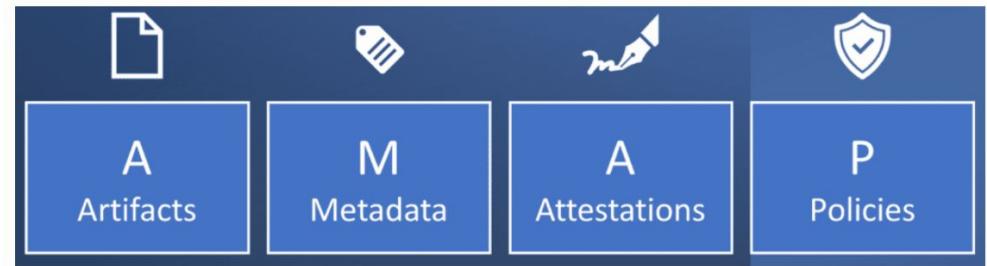
- Who guarantees that the key used for the signature is valid?
- Who guarantees that the guarantee we have for the key is valid?
- And so on... Turtles all the way down!



Supply Chain Security

Luckily, we can borrow ideas from “**classic**” DevSecOps to ensure the integrity of our model artefacts!

- Large landscape of projects!
- Sign **Artefacts**, **Metadata** and **Attestations**
- Apply **policies** to verify these



A MAP for software supply chain security

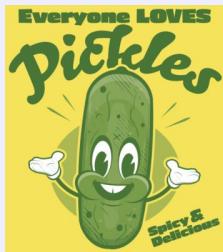
<https://www.cncf.io/blog/2022/04/12/a-map-for-kubernetes-supply-chain-security/>

Supply Chain Security

Artefacts

Our ML “binary” artefacts

- Pickles, ONNX, TF’s SavedModel, etc.



Metadata

- **Provenance Data:** Who trained this model? When did they train? What training pipeline did they use?
- **Software Bill of Materials (SBOMs):** What dataset was used? What package dependencies does the model have? We can use [CycloneDX ML-BOM](#) or (soon) [SPDX AI](#).
- **Vulnerability Scan Reports:** Known vulnerabilities within the model or its dependencies. We can use [VEX](#) format.

Attestations

Assurance that the artefacts and metadata come from the right persons and **haven’t been tampered with**.

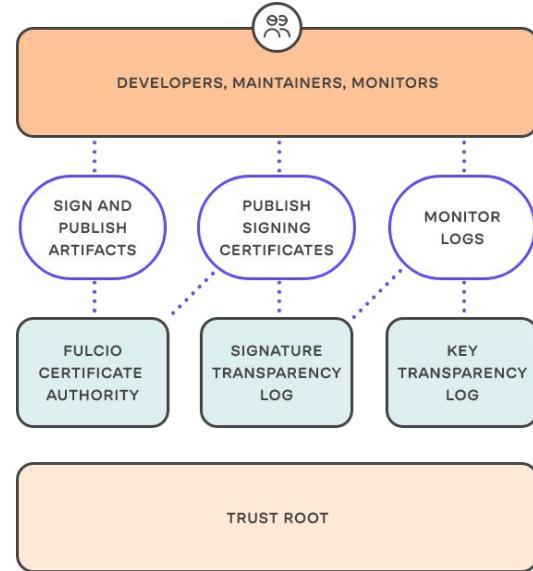
- Signature of Docker image
- **Signature** of our ML artefacts
- etc.



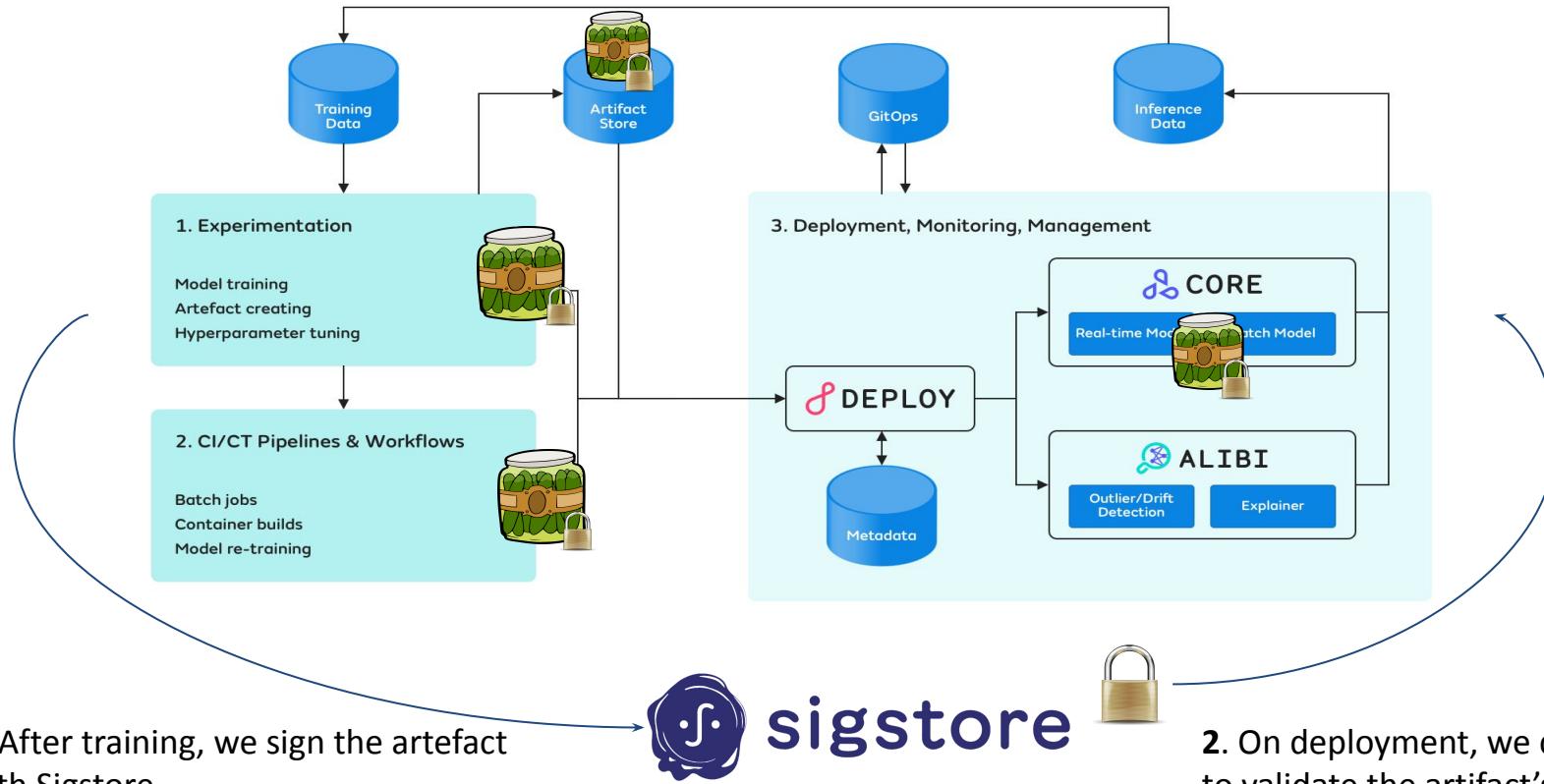
Supply Chain Security

Sigstore orchestrates the whole signing and verification process across a number of components:

- **Fulcio**: A free certificate authority.
- **Rekor**: Ledger with signed metadata to validate details.
- **OIDC Gateway**: Sigstore offers a simple OIDC Gateway, which can federate against Google, GitHub and Microsoft.
- The sigstore project runs hosted versions of the above, however they can also be run on-prem.



Tamper-proof Pickle jars



Demo



<https://github.com/adriangonz/sigstore-talk>

Next Steps

There don't seem to be many "*best practices*" around supply chain security for MLSecOps.

- No vendors provide signing - or the ability to apply policies (**yet**).
- Can we leverage other **OSS projects to apply policies** automatically?
- There's still work to do!

Monthly MLSecOps Wg Meetings

Join Monthly meetings:

forms.gle/N96w6h4wLDkfqsEs8

Key milestones

- Publish MLSecOps Top 10
- ML Security Standards
- Best Practices and Examples of Secure ML

Searched for "MLSecOps Working Group" in the search bar.

Calendars Projects Foundation Create ... Search

Pages / ... / Trusted AI Committee

MLSecOps Working Group

Created by Alejandro Saucedo just a moment ago

Attending Meetings

The current meetings take place monthly, alternating with the Trusted AI Committee meetings.

The meetings take place on the second Thursday of the month.

Please contact Alejandro Saucedo [@ethical.institute](#) if you would like to join the meetings and contribute to the MLSecOps Working Group.

You can also join the slack: LF AI & Data Foundation - [lfaifoundation.slack.com](#) - #mlsecops channel (coming soon) if you would like to engage in discussion about the group.

Working Group Members

For any required communication please contact the chairperson members:

- [@Alejandro Saucedo](#) (Working Group Chair) - Chief Scientist, The Institute for Ethical AI [@ethical.institute](#)
- [@Animesh Singh](#) (Trusted AI Committee Chair) - Chief Technology Officer, IBM [singhan@us.ibm.com](#)

Below are a list of core working group members - all affiliations listed for identification purposes only:

First Name	Last Name	Company	Position	Email Address
Beat	Buesser	IBM Research	Research Staff Member	beat.buesser@ie.ibm.com
Matthieu	Maitre	Microsoft	Principal Software Developer	mmaitre@microsoft.com
Aankur	Bhatia	IBM	Chief Data Scientist, IBM Security	abhatia@us.ibm.com
Karen	Bennet	Quality Craft	VP	bennetk@yahoo.com

Thank You

Adrian Gonzalez-Martin

Head of ML Serving

agm@seldon.io

