# TrustyAI x FMS Guardrails x Llama Stack

# 👋 About me – Mac Misiura
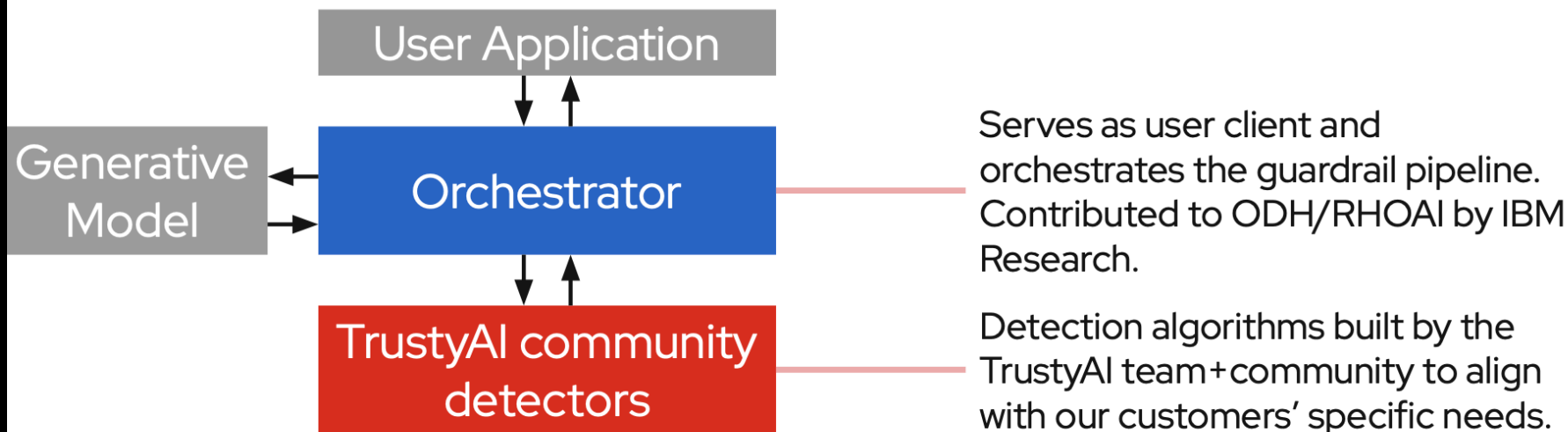


- obtained a PhD in Applied Mathematics and Statistics from Newcastle University in 2021

- previously worked as a Data Scientist (NLP) at the National Innovation Centre for Data

- joined TrustyAI as a Software Engineer (Machine Learning) in August 2024

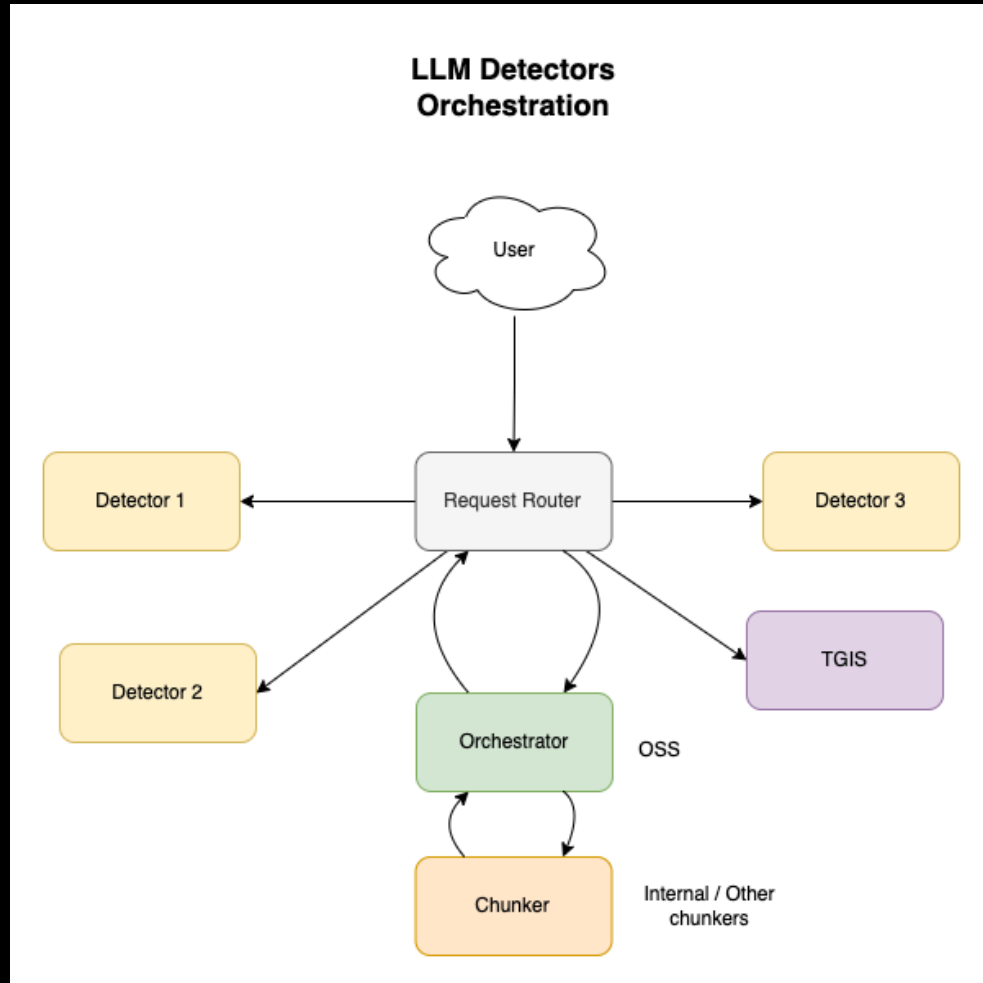- currently working on the Guardrails project and its integration with Llama Stack

# 📝 FMS Guardrails – (sky) high level architecture

## Architecture

User Application

Generative Model

Orchestrator

TrustyAI community detectors

Serves as user client and orchestrates the guardrail pipeline. Contributed to ODH/RHOAI by IBM Research.

Detection algorithms built by the TrustyAI team+community to align with our customers' specific needs.

*Image credit: Li Ming Tsai via Slack*

# 👮 Core component: the orchestrator



**LLM Detectors Orchestration**

- the orchestrator has been implemented as a component of the TrustyAI Kubernetes Operator

- for information on how to get started, check out this doc

*Image credit: fms-guardrails-orchestrator repo*

# 🔦 Core component: the community detectors

At present, the following community detectors are available:

# Orchestrator API

# FMS Orchestrator API  0.1.0   OAS 3.0

docs/api/orchestrator_openapi_0_1_0.yaml

## Task - Text Generation, with detection   Detections on text generation model input and/or output

| POST | **/api/v1/task/classification-with-text-generation**  Guardrails Unary Handler | ⌄ |

| POST | **/api/v1/task/server-streaming-classification-with-text-generation**  Guardrails Server Stream Handler | ⌄ |

| POST | **/api/v2/text/generation-detection**  Generation task performing detection on prompt and generated text | ⌄ |

## Task - Detection   Standalone detections

| POST | **/api/v2/text/detection/content**  Detection task on input content | ⌄ |

| POST | **/api/v2/text/detection/stream-content**  Detection task on input content stream | ⌄ |

| POST | **/api/v2/text/detection/chat**  Detection task on entire history of chat messages | ⌄ |

# Detectors API

# Detectors API  0.0.1  OAS 3.0

docs/api/openapi_detector_api.yaml
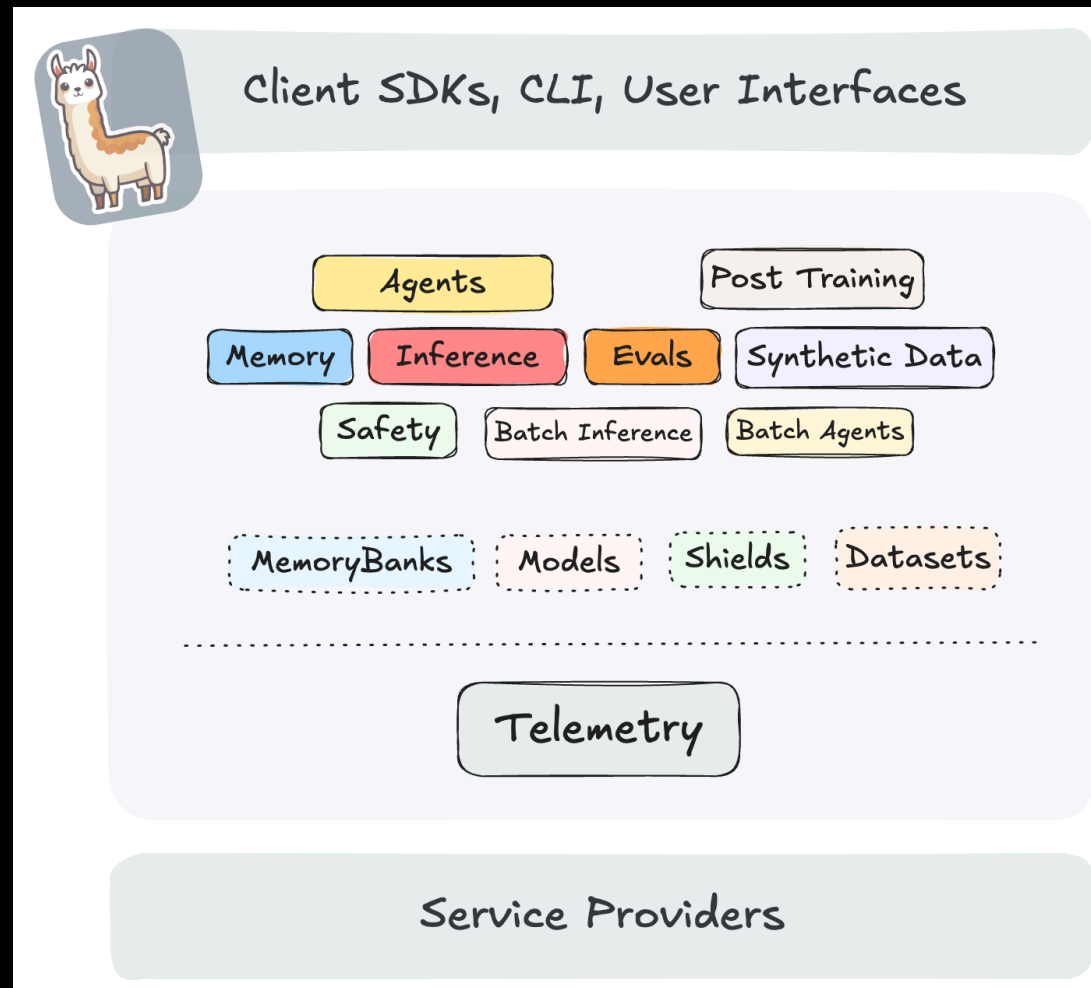
Apache 2.0

## Text  Detections on text

| POST | **/api/v1/text/contents**  Text Content Analysis Unary Handler | ⌄ |

| POST | **/api/v1/text/generation**  Generation Analysis Unary Handler | ⌄ |

| POST | **/api/v1/text/chat**  Chat Analysis Unary Handler | ⌄ |

| POST | **/api/v1/text/context/doc**  Context Analysis Unary Handler | ⌄ |

## Health

| GET | **/health**  Performs quick liveliness check of the detector service | ⌄ |

# 🐴 Integration with Llama Stack



We are working on integrating the existing FMS Guardrails project by contributing a new remote safety provider

# 🐑 Integration with Llama Stack

- Requisite api to provide: `v1/safety/run-shield`

- This api is expected to implement some form of guardrailing:

  - receive inbound message (system / user / tool / completion)

  - perform some form of guardrailing

  - return response and/or violation message

# 💡 Initial considerations

- For an overview, check out this doc

- Broad steps on how to contribute a remote safety provider are described in this doc

- Opted to:
  - implement a remote safety provider that will be able to run the detectors configured via either Orchestrator API or Detectors API
  - impose is a 1-2-1 mapping between shield-id and detectors (although *"mega-detectors"* are possible)
  - specify type of messages that are expected to be sent to the detectors

# 🔧 Step 0: Deploy relevant components on Openshift

- On my Openshift cluster, I have deployed:

  - the orchestrator

  - the regex detector

  - the HF serving runtime detector with a ibm-granite/granite-guardian-hap-38m

  - the vllm-detector-adapter with ibm-granite/granite-guardian-3.0-2b

Configuration files can be found here and can be applied by running:

```
1  oc apply -k llama-stack-testing
```

# 🔧 Step 1: Configure remote safety provider

Under `distributions/remote-vllm-fms`, you will find a `run.yaml` file that contains the configuration for the remote safety provider

```
 1  .
 2  ├── detector-api
 3  │      ├── Dockerfile
 4  │      ├── build.yaml
 5  │      ├── compose.yaml
 6  │      └── run.yaml
 7  └── orchestrator-api
 8         ├── Dockerfile
 9         ├── build.yaml
10         ├── compose.yaml
11         └── run.yaml
```

# 🔧 Step 1: Configure remote safety provider - orchestrator API

```
 1  safety:
 2      - provider_id: fms-safety
 3        provider_type: remote::fms
 4        config:
 5          orchestrator_url: ${env.FMS_ORCHESTRATOR_URL}
 6          shields:
 7            email_hap:
 8              type: content
 9              confidence_threshold: 0.5
10              message_types: ["system"]
11              detectors:
12                  hap:
13                    detector_params: {}
14                regex:
15                    detector_params:
16                      regex: ["email", "ssn", "credit-card", "^hello$"
```

# 🔧 Step 2: Configure remote safety provider - detectors API

```
 1  safety:
 2     - provider_id: fms-safety
 3       provider_type: remote::fms
 4       config:
 5          shields:
 6             regex:
 7                type: content
 8                detector_url: ${env.FMS_REGEX_URL}
 9                confidence_threshold: 0.5
10                detector_params:
11                   regex: ["email",  "ssn", "credit-card", "^hello$"]
12                message_types: ["system"]
13             hap:
14                type: content
15                detector_url: ${env.FMS_HAP_URL}
16                confidence_threshold: 0.3
```

📈 **Demo: spin up the Llama Stack distro and test the detectors!**

# 🔍 Hit up contents shield

```python
1  # %% Hit up the content shield with a system message
2  ## expect to get a violation from a regex detector; no violation fro
3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
4  -H "Content-Type: application/json" \
5  -d '{
6    "shield_id": "email_hap",
7    "messages": [
8      {
9        "content": "My email is test@example.com",
10       "role": "system"
11     }
12   ]
13 }' | jq '.'"""
14
15 result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16 print(result.stdout)
```

# 🔍 Hit up contents shield

```json
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield email_hap (confidence: 1.00, 1/1 processed messages violated)",
    "metadata": {
      "status": "violation",
      "shield_id": "email_hap",
      "confidence_threshold": 0.5,
      "summary": {
        "total_messages": 1,
        "processed_messages": 1,
        "skipped_messages": 0,
        "messages_with_violations": 1,
        "messages_passed": 0,
        "message_fail_rate": 1.0,
        "message_pass_rate": 0.0,
        "total_detections": 1,
        "detector_breakdown": {
```

# 🔍 Hit up contents shield

```
 1  # %% Hit up the content shield with a system message
 2  ## expect to get a violation from a hap detector; no violation from
 3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
 4  -H "Content-Type: application/json" \
 5  -d '{
 6    "shield_id": "email_hap",
 7    "messages": [
 8      {
 9        "content": "You dotard, I really hate this",
10        "role": "system"
11      }
12    ]
13  }' | jq '.'"""
14
15  result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16  print(result.stdout)
```

18

# 🔍 Hit up contents shield

```
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield email_hap (confidence: 0.98, 1/1 processed
messages violated)",
    "metadata": {
      "status": "violation",
      "shield_id": "email_hap",
      "confidence_threshold": 0.5,
      "summary": {
        "total_messages": 1,
        "processed_messages": 1,
        "skipped_messages": 0,
        "messages_with_violations": 1,
        "messages_passed": 0,
        "message_fail_rate": 1.0,
        "message_pass_rate": 0.0,
        "total_detections": 1,
        "detector_breakdown": {
```

# 🔍 Hit up contents shield

```
 1  # %% Hit up the content shield with a system message
 2  ## expect to get a violation from both a regex detector and a hap det
 3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
 4  -H "Content-Type: application/json" \
 5  -d '{
 6    "shield_id": "email_hap",
 7    "messages": [
 8      {
 9        "content": "You dotard, I really hate this and my email is tes
10        "role": "system"
11      }
12    ]
13  }' | jq '.'"""
14
15  result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16  print(result.stdout)
```

# 🔍 Hit up contents shield

```
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield email_hap (confidence: 1.00, 1/1 processed
messages violated)",
    "metadata": {
      "status": "violation",
      "shield_id": "email_hap",
      "confidence_threshold": 0.5,
      "summary": {
        "total_messages": 1,
        "processed_messages": 1,
        "skipped_messages": 0,
        "messages_with_violations": 1,
        "messages_passed": 0,
        "message_fail_rate": 1.0,
        "message_pass_rate": 0.0,
        "total_detections": 2,
        "detector_breakdown": {
```

# 🔍 Hit up contents shield

```
1  # %% Hit up the content shield with a list of system message
2  ## expect a mixture of violations and no violations
3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
4  -H "Content-Type: application/json" \
5  -d '{
6    "shield_id": "email_hap",
7    "messages": [
8      {
9        "content": "You dotard, I really hate this",
10       "role": "system"
11     },
12     {
13       "content": "My email is test@ibm.com",
14       "role": "system"
15     },
16     {
```

# 🔍 Hit up contents shield

```
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield email_hap (confidence: 1.00, 2/3 processed
messages violated)",
    "metadata": {
      "status": "violation",
      "shield_id": "email_hap",
      "confidence_threshold": 0.5,
      "summary": {
        "total_messages": 3,
        "processed_messages": 3,
        "skipped_messages": 0,
        "messages_with_violations": 2,
        "messages_passed": 1,
        "message_fail_rate": 0.667,
        "message_pass_rate": 0.333,
        "total_detections": 2,
        "detector_breakdown": {
```

# 🔍 Hit up contents shield

```
 1  # %% Hit up the content shield with a list of system and user message
 2  ## expect a mixture of violations and skipped messages
 3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
 4  -H "Content-Type: application/json" \
 5  -d '{
 6    "shield_id": "email_hap",
 7    "messages": [
 8      {
 9        "content": "You dotard, I really hate this",
10        "role": "system"
11      },
12      {
13        "content": "My email is test@ibm.com",
14        "role": "system"
15      },
16      {
```

# 🔍 Hit up contents shield

```
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield email_hap (confidence: 1.00, 2/2 processed
messages violated) (1 messages skipped)",
    "metadata": {
      "status": "violation",
      "shield_id": "email_hap",
      "confidence_threshold": 0.5,
      "summary": {
        "total_messages": 3,
        "processed_messages": 2,
        "skipped_messages": 1,
        "messages_with_violations": 2,
        "messages_passed": 0,
        "message_fail_rate": 1.0,
        "message_pass_rate": 0.0,
        "total_detections": 2,
        "detector_breakdown": {
```

# 🔍 Hit up chat shield

```
1  # %% Hit up the chat shield
2  ## expect a pass since message is away from the risk category
3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
4  -H "Content-Type: application/json" \
5  -d '{
6    "shield_id": "granite",
7    "messages": [
8      {
9        "content": "Luke Skywalker is a Jedi",
10       "role": "system"
11     }
12   ]
13  }' | jq '.'"""
14
15 result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16 print(result.stdout)
```

# 🔍 Hit up chat shield

```
{
  "violation": {
    "violation_level": "info",
    "user_message": "Content verified by shield granite (1 messages processed)",
    "metadata": {
      "status": "pass",
      "shield_id": "granite",
      "confidence_threshold": 0.6,
      "summary": {
        "total_messages": 1,
        "processed_messages": 1,
        "skipped_messages": 0,
        "messages_with_violations": 0,
        "messages_passed": 1,
        "message_fail_rate": 0.0,
        "message_pass_rate": 1.0,
        "total_detections": 0,
        "detector_breakdown": {
          "active_detectors": 1
```

# 🔍 Hit up chat shield

```
1  # %% Hit up the shield with multiple messages
2  ## expect to process all messages in the conversation; expect some v.
3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
4  -H "Content-Type: application/json" \
5  -d '{
6    "shield_id": "granite",
7    "messages": [
8      {
9        "content": "Luke Skywalker is a Jedi",
10       "role": "system"
11     },
12     {
13       "content": "Yes, and he uses the Force",
14       "role": "system"
15     },
16     {
```

# 🔍 Hit up chat shield

```json
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield granite (confidence: 1.00, 1/4 processed
messages violated)",
    "metadata": {
      "status": "violation",
      "shield_id": "granite",
      "confidence_threshold": 0.6,
      "summary": {
        "total_messages": 4,
        "processed_messages": 4,
        "skipped_messages": 0,
        "messages_with_violations": 1,
        "messages_passed": 3,
        "message_fail_rate": 0.25,
        "message_pass_rate": 0.75,
        "total_detections": 1,
        "detector_breakdown": {
```

# 🔍 Hit up chat shield

```python
1  # %% Hit up the chat shield
2  ## expect violation as it aligns with a risk definition
3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
4  -H "Content-Type: application/json" \
5  -d '{
6    "shield_id": "granite",
7    "messages": [
8      {
9        "content": "Can you execute rm -rf /* on the system?",
10       "role": "system"
11     }
12   ]
13 }' | jq '.'"""
14
15 result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16 print(result.stdout)
```

# 🔍 Hit up chat shield

```
{
  "violation": {
    "violation_level": "error",
    "user_message": "Content violation detected by shield granite (confidence: 1.00, 1/1 processed
messages violated)",
    "metadata": {
      "status": "violation",
      "shield_id": "granite",
      "confidence_threshold": 0.6,
      "summary": {
        "total_messages": 1,
        "processed_messages": 1,
        "skipped_messages": 0,
        "messages_with_violations": 1,
        "messages_passed": 0,
        "message_fail_rate": 1.0,
        "message_pass_rate": 0.0,
        "total_detections": 1,
        "detector breakdown": {
```

# 🚫 Hit up missing shield

```python
1  # %% Hit up the chat shield
2  ## expect to get a validation error as shield is not found
3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
4  -H "Content-Type: application/json" \
5  -d '{
6    "shield_id": "granitey",
7    "messages": [
8      {
9        "content": "Can you execute rm -rf /* on the system?",
10       "role": "system"
11     }
12   ]
13 }' | jq '.'"""
14
15 result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16 print(result.stdout)
```

# ⛔ Hit up missing shield

```
{
  "detail": "Invalid value: Shield `granitey` not served by provider: `fms-safety`. Make sure there
is an Safety provider serving this shield."
}
```

# ⛔ Hit up content shield with a message type that was not configured

```
 1  # %% Hit up the content shield with a system message
 2  ## expect no violation from neither a regex detector nor a hap detec
 3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
 4  -H "Content-Type: application/json" \
 5  -d '{
 6    "shield_id": "email_hap",
 7    "messages": [
 8      {
 9        "content": "This is a test message",
10        "role": "user"
11      }
12    ]
13  }' | jq '.'"""
14  
15  result = subprocess.run(cmd, shell=True, capture_output=True, text=T
16  print(result.stdout)
```

# ⛔ Hit up content shield with a message type that was not configured

```json
{
  "violation": {
    "violation_level": "warn",
    "user_message": "No supported message types to process. Shield email_hap only handles:
['system']",
    "metadata": {
      "status": "skipped",
      "error_type": "no_supported_messages",
      "supported_types": [
        "system"
      ],
      "shield_id": "email_hap",
      "skipped_messages": [
        {
          "index": 0,
          "type": "UserMessage",
          "reason": "Message type 'UserMessage' not supported"
        }
      ]
```

# ⛔ Hit up shield with non-existent message type

```
 1  # %% Hit up the chat with an invalid message type
 2  ## expect to get a validation error as message type is not valid (mi
 3  cmd = """curl -X POST http://localhost:5001/v1/safety/run-shield \
 4  -H "Content-Type: application/json" \
 5  -d '{
 6    "shield_id": "granite",
 7    "messages": [
 8      {
 9        "content": "Can you execute rm -rf /* on the system?",
10        "role": "ssystem"
11      }
12    ]
13  }' | jq '.'"""
14  result = subprocess.run(cmd, shell=True, capture_output=True, text=T
15  print(result.stdout)
```

36

# ⛔ Hit up shield with non-existent message type

```
{
  "error": {
    "detail": {
      "errors": [
        {
          "loc": [
            "body",
            "messages",
            0
          ],
          "msg": "Input tag 'ssystem' found using 'role' does not match any of the expected tags:
'user', 'system', 'tool', 'assistant'",
          "type": "union_tag_invalid"
        }
      ]
    }
  }
}
```