# FAM Moderation Solution Discussion

## Background

FAM is leveraging EverC and Incopro to do AUP monitoring and BwP branding on the site monitoring. However, the existing solutions are inadequate for our long-term needs, primarily due to limited functionality, lack of data transparency and the escalating costs associated with relying on third-party services. Also we have identified a P0 requirement which AUP/Incopro solution wouldn't support, to help mitigate the risk for counterfeit listings in FOLEX from 12/15/23. On March 2023, FAM team created PR/FAQ to build flexible monitoring and moderation capabilities to support use-cases ranging from policy violation detection to site content/structure optimization. This involves implementing two distinct systems: a dedicated monitoring system to extract structured information from a website at pre-defined cadence or post a trigger activity; and a separate moderation system to handle evaluate data against a set of policies. For more background, please refer to doc Monitoring and Moderation High Level Discussion.

In this quip, our primary focus will be on the moderation aspect. The moderation system will have the responsibility to evaluate compliance of extracted object/entity against specific defined policies. The goal for this doc is to provide clarity regarding the requirements and explore various potential options, along with their respective advantages and disadvantages. We will make an informed decision based on reviews and keep going to focus on next steps with high level designs.

## Requirements

### END FEATURE PRIORITY

The FAM moderation system will build the following features with priority from high to low:

1. **[Feature 1] Moderate Santos catalog data for LimaPDP and Hosted Widget Post Order Pages:** Santos catalog team will send info to RiskEvaluator with the catalog data, which includes Lima PDP and hosted Widget Post Order Pages like BwP Checkout Page, Order Confirmation Page, Order Tracking Page. Moderation needs to detect any non-compliant info in catalog.
2. **[Feature 2] Moderate FOLEX listings for Automated Brand Protection and Brand Gating controls:** If content (image or text) that is a trademark or brand logo of a well known brand. Moderation needs to check these results with list of allow listed brands, then human investigator will decide if it is a counterfeit brand. (Extra efforts are needed because current existed ML models didn't cover this use case)
3. **[Feature 3] Moderate merchant declared URLs across** all business products**:** All merchants onboarded with BwP products need to comply BwP AUP, so after continuously monitoring merchant websites, moderation needs to detect any non-compliant contents.

### BUSINESS REQUIREMENTS

1. [P0] As a moderation system user, I can moderate content in the form of texts and images for product attributes level.
2. [P0] As a moderation system user, I can moderate text content for top 20 highest risk brands in FOLEX.
3. [P0] As a moderation system user, I can utilize default policies to trigger the moderation function.
4. [P0] As a moderation system user, I can receive moderation results asynchronously, with confidence scores and other analyzes output from the moderation model and an outcome label as APPROVE, REJECT, or MANUAL_REVIEW.
5. [P0] As a user who initializes a moderation function, I should be able to configure confidence score thresholds to determine the moderation results.
   a. If the confidence score of one content is above the threshold, the outcome should be labeled as the result from moderation model: APPROVE means the content adheres to the moderation guidelines; REJECT indicates that it is not qualified.

      b. Conversely, if a content achieves a confidence score lower than the threshold, with ambiguous/uncertain, the result should be labeled as MANUAL_REVIEW.

6. [P0] As a user of the moderation system, I can <u>manually investigate in Paragon</u> whenever the moderation model generates MANUAL_REVIEW results. The paragon case includes corresponding confidence scores and an analysis of the output produced by the moderation model. Each moderation request will only generate one paragon case.  For example, the moderation request for a specific domain or merchant, we will generate one paragon case with a list of URLs which marked as MANUAL_REVIEW.

7. [P1] As a moderation system user, I can select <u>multiple policies</u> from a predefined list of BwP policies.

8. [P1] As a user who initializes a moderation function, I can configure the elements displayed in the paragon case, including the ability to change the title and <u>customize paragon contents</u> to align with my use case.

9. [P1] As a FAM moderation investigators, I can assign approval or rejection when resolving paragon cases. My investigation <u>decisions will send back</u> to the system saved as the final moderation results.

10. [P1] As a user who initializes a moderation function,  I can use a new policy with an allowlist, denylist and manual review list to approve or reject content based on specific keywords, as a straightforward content moderation based on keyword criteria.

11. [P2] As a moderation system user, I can moderate content in the form of <u>videos</u> with URLs.

12. [P2] As a moderation system user, I can moderate brand logo/trademark<u> image for FOLEX counterfeit brands.</u>

13. [P2] As a FAM moderation investigators, I should feed case results back to ML training so that the ML models can learn and get better. (Depends on performance in P0)

14. [NEW called out by Ryan][P1] Feed manual moderation results to the system to avoid duplicate moderation results.

**SYSTEM REQUIREMENTS**

1. Moderation system has inputs as plain text, image/video URLs, and the S3 folder path to find a list of files. All contents are <u>already formatted</u> according to the predefined schema established by FAM.

2. Moderation system is developed as a new risk evaluation function, which leverages FAM Risk Evaluator to link with the Santos domain team, allowing them to trigger evaluations and take appropriate actions based on the RiskOutcome output within moderation function result.

3. Moderation system supports <u>various moderation methods</u>, like regex expression check, machine learning models related to natural language processing (NLP) and image recognition (IR) algorithms. Additionally, the moderation system is designed to be extensible,  allowing the easy addition of new models.

4. The moderation system provides support for both asynchronous[p0] and synchronous[p1] moderation review processes.

      a. If the models can generate results within x seconds, we will immediately return the results to the client.

      b. If the processing time longer than x seconds, we will return the results as IN_REVIEW and send the moderation review result asynchronously.

      c. If manual review is enabled, the time from manual review will determine the final outcome and completion of the moderation process.

# Design Decisions Approaches

We listed current solutions and few approaches we considered for moderation.

| | Current Solution | Approach 1 | Approach 2 | Approach 3 | Approach 4 |
|---|---|---|---|---|---|
| Name | Apay/EverC | Reuse Amazon Classification and Policy Platform(CPP) | Create FAM moderation platform | Create FAM platform with Oberon ML Model | Reuse Oberon Platform |
| Description | FAM is using EverC and Incopro to do AUP monitoring checking | Reuse CPP platform which provides self-service classification of products in the catalog and management of policies. | Build own platform using existing fundamental methods and create FAM owned ML models. | Build FAM-owned platform with Oberon models from Review team. | Oberon moderation system is the service Review team build to serve moderation requirements. |
| Moderation methods | Unknown | FAM can reuse existed classifications If it's covered in our current use case already | FAM will build ML models. Needs collaboration with ML scientists to develop models based on our use cases from scratch | Oberon ML model is already onboarded 5 machine learning models from CS team and we can reuse their sagemaker endpoint to use their ML models | Platform provides Oberon ML model as a set of comprehensive moderation policies proven in Amazon's decades of e-commerce operation |
| Manual Moderate | Once EverC send the violations back, FAM team will review to determine if the violation is true | CPP has the UI portal for investigators. Need extra efforts to connect CPP platform with paragon | Needs to connect paragon for manual moderation during implementation | Needs to connect paragon for manual moderation during implementation | Oberon Platform has the UI for investigators. Extra efforts needed to connect paragon case |
| Flexibility | **N/A**<br><br>No visibility and flexibility on what pages and the content they scanned | **+**<br><br>New feature requests need to fill Project Intake Process template | **++++**<br><br>FAM can respond rapidly and take action items with new features internally | **+++**<br><br>FAM can respond rapidly and take action items with new features internally. But ML models maintenance needs more help from Review team. | **++**<br><br>FAM needs Review team to support new features in Oberon service. |
| Development Time estimation | N/A | **+++**<br><br>CPP has a detailed onboarding plan to support new project. On-demand classification needs extra SDE efforts. | **+++++**<br>1. Code implementation<br>2. Full Security review<br>3. Large size efforts from Scientists | **+++**<br>1. Code implementation<br>2. Full Security review<br>3. Medium size efforts from Scientists in Review team | **++**<br>1. SDEs efforts from Review teams<br>2. Delta Security review<br>3. Medium size efforts from Scientists in Review team |
| Ownership and Maintenance | N/A | FAM will be a user in CPP. New feature addition/enhancement request in CPP platform needs to follow intake process | FAM owns and maintains platform and ML models | FAM owns and maintains the Platform; Review team owns ML model and takes over new model requests from FAM | Review team owns platform and ML model. FAM call Oberon as a dependency |
| Long term plan | Cannot fit for [Feature 2] | It's not a long term solution | Can be extended to a santos wide moderation system | Can be extended to a santos wide moderation system | Can be extended to a santos wide moderation system |
| Pros | 1. + EverC is the third-party which strengthen and streamline risk management<br>2. + Friendly UI with moderation results. | 1. + Well maintained platform with clear UI which is widely used in Amazon.<br>2. + Move fast: FAM can reuse existed classification from offensive team to get moderation results without too much SDE/ML efforts | 1. + FAM will have complete ownership for the moderation platform<br>2. + FAM has high flexibility to customize ML models and add new features in the future | 1. + FAM will have complete ownership for the moderation system architecture<br>2. + Move fast: Leveraging existing ML models from review teams<br>3. + Accurate: Beta version implementation duplicates Amazon CS ML model with 100% decision match rate. | 1. + Move fast: Oberon has already used Rekognition and ML models to moderate which mentioned in approach 2<br>2. + Accurate: Beta version implementation duplicates Amazon CS ML model with 100% decision match rate.<br>3. + Avoid duplicate work to build moderation system. |
| Cons | 1. - Lack of data transparency<br>2. - Limited of functionality<br>3. - Ballooning Cost | 1. - CPP supports data in Catalog with ASIN based classifications only. Catalog data cannot cover header/footer and non-product pages in [Feature 3]<br>2. - Uncertainty of the timeline within FAM's new feature requests<br>3. - Extra costs and efforts associated when integrating with Amazon CDO platform with Santos data. | 1. - Substantial efforts from SDEs due to the complex nature of integrating various methods for text, image, and video moderation together<br>2. - Substantial efforts from machine learning scientists to create ML models from scratch | 1. - Substantial efforts from FAM SDEs to build platform<br>2. - Scientists from Review team need to provide supports in maintaining moderation accuracy for FAM. | 1. - Extra away team efforts to extend inputs limited, and support paragon cases in domain level<br>2. - The Review team owns the core moderation logic. For any new features, external team support will be required. |
| Decision | Goal is to replace by FAM moderation service | Not recommended for long term | Not recommended | Considering | Accept |

We only listed the key decision-making elements here. Details for each approach are shared in section Approaches details. For additional factors such as processing time, cost, and policy supports, you can find further information in Approaches comparing - Extend.

## Summary

No matter FAM creates our own ML models/platform, or reuse Review's models or platform in approach 3/4, we are still targeting to build a santos-wide moderation system from long-term respective. This santos moderation system should be used by other Santos teams and it's centralized with multiple moderation functions like RegEx, ML, keyword-based rules, etc. In order to move

fast in the initial phase and take minor efforts to reach the long-term goal, FAM will choose to adopt Approach 4 as the starting point for the moderation process.

**IMPLEMENTATION PHASES**

Here are the implementation phases for integrating the Oberon service into FAM:
**[Phase 1] - One time execution with ML model - 08/17 for POC**
Oberon serves as a well-trained moderation service for reviews data; however, it remains unknown to moderate data from Santos catalog. At an early stage of the implementation, we will initiate a POC to moderate catalog data within existing ML model. The goal is to validate its accuracy and suitability for our requirements, and also to give some feedbacks to scientists to train the model.

- Input: catalog data files provided by FAM (Full catalog data)
- Output: Preliminary moderation results generated by the ML model

**[Phase 2] - Periodically execution with santos catalog data**
FAM is trying to moderate Santos catalog data for LimaPDP and Hosted Widget Post Order Pages in a timely manner. So we aim to launch a simple version of the moderation system to periodically moderate Santos catalog data. With weekly or daily moderation, FAM will be able to initiate some investigations for non-compliant case, thereby reducing potential risks to protect BwP reputation. Additionally, as we progress through Phase 3, the increasing number of data points from the Santos catalog will provide valuable training data to further enhance the ML model.

- Input: catalog data files provided by FAM
- Output: Oberon moderation reports (eg. csv) for FAM investigator

**[Phase 3] - Integrate with Risk Evaluator workflow and call Oberon with contents through API <end of 2023>**
In the final phase of implementation, FAM will utilize API calls to retrieve moderation results from Oberon with FAM Risk Evaluator workflow. Oberon will serve as the core moderation system to deliver moderation results.

- Input: API call with contents which need to moderate
- Output: results from Oberon moderation system

The implementation phases above is the draft version after we finalize the approaches. For the next step, FAM will engage Review team to collaboratively determine the most suitable ways to execute each phase in Oberon and FAM. More detailed plan will be the major part in the High-Level Design.

**ACTION ITEMS**

1. [Done] Set up the meeting with Thibault, and bring the solutions to FAM. -Xiaoxi Bai
2. [ETA Pending] Talk with offensive team to get the keywords list to classify. - Aniruddh Menon
3. [Done] Sync with review team for moderation ML model. - Xiaoxi Bai
   a. Peng - Oberon: how many efforts to add/support new ML model.
   b. Counterfeit needs new ML model to check brand name
4. [ETA 08/03] Sync with catalog team for moderation ML model. - Xiaoxi Bai
   a. Clarify FAM's requirement
5. [ETA 08/11] Start the System Design and clarify following topics - Xiaoxi Bai
   a. ~~one time moderation within RE function, or can with ml model only.(backfill or manual triggered execution)~~
   b. build paragon inside or outside evaluation function
   c. the way to validate merchant can or cannot sell this brand offcially based on soucelist.(regard unauthorized merchant as counterfeit)
   d. Sizing for approach3/approach4; sizing put

  e. SLA, securtity

  f. how we handle counterfeit

  g. We might want to add a milestone to support FOLEX data as well, but can be tracked as a future deliverable after knowing better on timelines

6. [ETA 08/18] Review with FAM team and review team.

---

(We can stop reading here. Rest of doc can be the offline reading if you are interested)

# Approaches details

## APPROACH 1: REUSE AMAZON CLASSIFICATION AND POLICY PLATFORM(CPP)

The Classification and Policy Platform (CPP) provides self-service classification of products in the catalog and management of policies based on these classifications by business users. It has an audit platform with multiple classification techniques including ML, keyword-based rules, similarity algorithms, image classification combined with human classifications on product data, which means we have the option to onboard ML models for product classification and train them ourselves. Also, if there are existing classifications in the system that cater to FAM's use cases, we can leverage them without too much SDE/ML scientists efforts to create new classifications from scratch.

Although Approach 1 can quickly deliver results for short-term based on Santos product Asins, it is not advisable for long-term use in FAM due to:
1) Its limited support for data in the catalog only. FAM is to ensure comprehensive moderation of Santos merchant sites. We are considering merchant egregious content can be present not only in product titles/details, but also in various other parts of the content like headers, footers and other non product detail pages etc.
2) In order to process Santos catalog data and store results in CPP platform, exceptions from Santos data handling team are required. Synced with them in the Office hour that, it's challenging because santos promise to merchants on not sharing their data with Amazon wide.
3) CPP's 2023 roadmap has already been established, and they cannot accommodate new requests without clear impact and severity quantification. Therefore, it is vital to take into account the uncertainties and potential delays associated with submitting requests to support additional Santos contents, or new features at a later stage. Given the costs and efforts(security approval exceptions) required for integrating with Amazon CDO, the recommended approach is to build service without dependencies from Amazon CDO services. This allows for greater flexibility while minimizing potential obstacles.

See CPP Demo
here: https://amazon.awsapps.com/workdocs/index.html#/document/52bb22597b15de9f2454483af47be2b2c1be1408527a30072b56ee4da06b948e

**Key points: CPP is insufficient to cover all contents from Santos merchant sites in our feature3. Also, we need extra data handling exceptions and security efforts to engage santos catalog data into CDO services.**

## APPROACH 2: BUILD FAM PLATFORM WITH RULE BASED AND ML MODERATION MODULE

Approach 2 requires FAM to construct its own platform using existing fundamental methods to moderate contexts. We need to integrate various methods to handle different types of content, including text, image, and video. Some potential options for moderation methods are summarized in Moderation methods in approach 2.

FAM will have complete ownership and responsibility for the moderation service in approach 2. We have high flexibility to add

new features in the future. But it needs substantial efforts from SDEs in FAM due to the complex nature of integrating various methods for text, image, and video moderation. Also, it's requiring collaboration with machine learning scientists to develop models based on our use cases from scratch. To maximize reusability and efficiency in our implementation, We gave approaches 3 and 4, which focus on exploring opportunities to leverage existing ML models and functions from Santos-wide sources.

**Key points: FAM needs substantial efforts from engineers and dedicated scientists to build the whole moderation system from scratch.**

## APPROACH 3: BUILD FAM PLATFORM WITH MODELS FROM OBERON PLATFORM

Currently, Oberon used machine learning models from CS team which has been widely used for content moderation. They onboarded 5 machine learning models from CS team. For text moderation workflow, there are ReviewsTextModel which is used to predict APPROVE/REJECT decision for given text content, and ReviewsRejectReasonModel which used to predict the reason code for rejected review. And For image moderation, there are ReviewsNudityModel which focuses on detecting nudity in review image; ReviewsGlobalModel which detects inappropriate factors other than nudity; and ReviewOffTopicModel which compares the review image with product image and predict the relevancy. Here is the supported policies summary: Supported Policies in summary: 📄 Oberon Moderation Service Micro-Policy Support

Approach 3 makes up for the shortcomings from approach 2. Oberon beta version implementation duplicates Amazon CS ML model with 100% decision match rate. FAM can accelerate its progress by leveraging existing ML models. On the downside, some extra efforts are needed in FAM to build platform. Also, scientists from Review team need to provide supports for FAM in maintaining moderation accuracy, allowing the FAM team to fully concentrate on system architecture and implementation.

**Key points: FAM can leverage existing ML models from review teams to minimize the efforts from scientists. We still need to create moderation platform.**

## APPROACH 4: REUSE OBERON PLATFORM FROM REVIEW PLATFORM

Oberon moderation system is the service review team build to serve moderation requirements. This service has already provided a set of comprehensive moderation policies proven in Amazon's decades of e-commerce operation, to help customers without proper moderation policies setup up front.

Oberon platform used Rekognition and ML models which we considered in Approach 2. Reusing will avoid the duplicate work in FAM moderations. The designed reviewed with Kyle and functionality works well with ML models. In Approach 4, some additional efforts are required as the current usage of Oberon moderation system is limited to reviews only, with inputs limited to text and images. To align with our specific use case, we need to extend Oberon moderation platform to support file paths for processing batched messages. Additionally, the existing setup of Manual review within the moderation review needs to be reconfigured to create a paragon case based on a single domain URL or merchant. On the downside, approach 4 indicates Review team will have complete ownership and responsibility for the moderation service. FAM will entirely depend on Review team to support features in Oberon service. Even for any new features in later stages, external team support will be required.

**Key points: FAM can leverage existing ML models, and avoid some duplicate works to build similar moderation system in approach 3.**

## APPROACH NOT IN CONSIDER:  REUSE BERLIN MODERATION FROM COMMUNITY SHOPPING TEAM

Berlin Moderation team preserves shoppers' trust by ensuring that the content uploaded to Amazon is appropriate and relevant to the shopping experience. We define the policies and guidelines around appropriateness, relevancy to the shopping experience, and the categorization of the content. When contributors intentionally violate Moderation guidelines, we enforce them by warning and/or banning them from using Community Features. Berlin Moderation is already a rule based moderation system with detailed on-boarding guidance. Approach 5 is not recommend because:

1. FAM requires customization and extensibility in order to cater to different moderation policies beyond the scope of

Amazon. As potential clients of FAM may belong to various domains other than Amazon retail, the existing filter rules library, which is heavily retail-focused and established by the CS team, <u>might not adequately cover use cases</u>.

2. For optimal functionality, FAM should separate itself from <u>Amazon CDO dependencies</u> and operate as an independent service. Taking into account the cost and effort associated with integrating with Amazon CDO, as well as the need for security approval to create exceptions, building our own service is the preferred solution to effectively fulfill this objective.

## APPROACHES COMPARING - EXTEND

| | Current Solution | Approach 1 | Approach 2 | Approach 3 | Approach 4 |
|---|---|---|---|---|---|
| Name | Apay/EverC | Reuse Amazon Classification and Policy Platform(CPP) | Create FAM platform | Create FAM platform with Oberon ML Model | Reuse Oberon Platform |
| Data Inputs | Share Merchant info with domain URL | Now it supports data in Catalog with ASIN based classifications only | Will support inputs as text, image, video and file path | Will support text, image, video and file path | Now it supports text and image |
| Customzied Policies | No | Yes<br><br>Policy Authors can define, author and manage policies against classified ASINs | Yes | Yes | Yes |
| Customized Threashold | No | Yes | Yes | Yes | Yes |
| Moderation processing SLA without manual review | **++++ (monitoring + moderation)**<br><br>1-3 days to get results from EverC | +++<br><br>~2 hrs | ++<br><br><10s depends on text and policies | ++<br><br><10s depends on text and policies | ++<br><br><10s depends on text and policies |
| Cost | $$$$$<br><br>Using EverC and Incopro will cost ~$1.5M daily in 2024 | $<br><br>Need to pay hardware and provisioning cost based on CPP IMR table.<br><br>* Estimation: $0.875 to process one million asins per hour = ~$20 daily. | $$<br><br>Major cost are coming from sagemaker and Rekognition<br><br>* Estimation: If we use two 32vCPU 128 GB sagemaker Instances($0.88/hr) and process 30K images in rekognition($1.27). Total cost = ~$27.98 | $$<br><br>Same as approach 2 | $$<br><br>Same as approach 2 in review team |

# Appendix

## BUSINESS METRICS REQUIREMENT

We will work with FAM TPM to clarify metrics:

1. Total # of requests for moderation coming from Catalog, or Monitoring(Domain Url).
2. Total # of functions pending for MANUAL_REVIEW in paragon.
3. Total # of APPROVE, REJECT after manual review
4. Top 5 reasons/policies for the most REJECT cases
5. Top 5 policies for the most MANUAL_REVIEW cases
6. Total time for moderation function executions; Average time we take for manual reviewing.
7. The confidence score(0-99) for each policy APPROVE, REJECT contents.

## MODERATION METHODS IN APPROACH 2

### Text - Rule based RegEx filter
Rule based moderation module evaluates text content against RegEx filter rule, and determines approve/reject/manual result based on filter rule configuration. We can reuse filter rules library from Amazon Community Shopping(CS) team, which contains around 39k individual filter rules, each of them has a RegularExpression(RegEx) and corresponding moderation decision if matched. These RegEx filter rules have been categorized into multiple lists based on its attribute(e.g. HateSpeech, Spam, etc.) and targeting language(e.g. EN, ES). The rule-based moderation module generates rapid results, which we can promptly return

to the client as the prescreening moderation. However, for the remaining ML moderation results, we will send them back to the client asynchronously.

### Text - ML - Amazon Comprehend

Amazon Comprehend is a natural language processing (NLP) service that uses machine learning to discover insights from text. This module can predict whether the content adheres to moderation guidelines.
**Limits:** This method does not support image, video, and human moderation. Its functionality is limited to extracting sentiment and entities from unstructured text. The labeling focuses on entities rather than identifying inappropriate content. Only output labels, need additional logic to deduce final moderation decision.

### Text - ML - CPP AutoML

AutoML is a self-service ML solution for business users to train product classification models with custom precision/recall requirements, without the need for applied scientists. It automates feature engineering, feature selection and hyper parameter tuning and utilizes an ensemble model based on Logistic Regression and Random Forest classifiers for making predictions. AutoML uses ensemble of Logistic Regression model and Neural Network Language Modeling (NNLM) based embedding based RF model.
**Limits:** Not support for image, video moderation. User needs to provide example ASINs (aka labels) of members and non-members of the class using AutoML UI to train the model. (TODO: check whether it used in amazon catalog with asin only.)
Ref: https://w.amazon.com/bin/view/SelectionClassification/CPP_AI/AutoML/

### Image - Rekognition

Amazon Rekognition can be used to detect content that is inappropriate, unwanted, or offensive. Rekognition moderation APIsare used in social media, broadcast media, advertising, and e-commerce situations to create a safer user experience, provide brand safety assurances to advertisers, and comply with local and global regulations.  **Limits:** No text moderation support; Cannot predict off-topic content by comparing feedback image with original product image. Only output labels, need additional logic to deduce final moderation decision.

### Image/Text - ML - SageMaker with data labeling

Amazon SageMaker is a fully managed service that provides every developer and data scientist with the ability to build, train, and deploy machine learning (ML) models quickly. We can deploy machine learning models at any scale and highly customizable for any machine learning needs.  For image detection, Amazon SageMaker image classification algorithm uses a modified version of the Mxnet image classification algorithm. For example, with SageMaker image classification algorithm, we can predicts image and reject feedback image if the nudity score is above rejection threshold.

### FAM MODERATION INVESTIGATION

FAM moderation Investigations covered topics for:

- Oberon Service Background
- CPP Background
- Options in approach 3/approach 4 to reuse the Oberon platform

---

# Agenda for internal review with Kevin, Anirrudh and Troy 7/24

1. Clarify the requirements
2. Approches discussion
    a. Agreements for final state in moderation in Santos

      b. Clarify direction and ownership.

           i. Approach 1: Painpoint - Catalog data based on ASIN

           ii. Approach 2: Painpoint - Duplicate work, extra ML

Updates:

1. Update counterfeit use case in the requirement
2. Replied on the comments

**Follow up: 7/27**

Goal:
Finalize the solution discussion the solutions which can move fast and also extenable for other use case. Clarify the next step for moderation.

1. The summary from CPP
   a. Pros: 1. friendly UI, self onboarding → move fast  2. provide some classifications and get some quick results. Offencies team.  offensive. Existed model with more FAM use case. violation
   b. Cons: 1. Asin only in catalog 2. in amazon catalog data in
   c. DH: our promise to merchants on not sharing their data with Amazon.
   d. Not a good start point: 1. Long term → No 2. Short term → costing; even On demand
   e.
2. Phase 0
   a. Benefit? 1. move fast have the moderation results 2.
   b. Periodically execute data with model

# Review with Thibault 7/31

1. Got the aligment with approach 4
2. Forwarded the review team sync meeting to Thibault

# Review with Thibault 8/1

Pavan: as far as I know there are 2 types of Lima PDP urls. One stored in Catalog, and one generated during runtime by X builder team. Do we also want to moderate other PDP url types?