# Looper: an end-to-end ML platform for product decisions

IGOR MARKOV, HANSON WANG, NITYA KASTURI, SHAUN SINGH, MIA GARRARD, YIN HUANG, SZE WAI YUEN, SARAH TRAN, ZEHUI WANG, IGOR GLOTOV, TANVI GUPTA, PENG CHEN, BOSHUANG HUANG, XIAOWEN XIE, MICHAEL BELKIN, SAL URYASEV, SAM HOWIE, EYTAN BAKSHY, and NORM ZHOU, Meta, USA

Modern software systems and products increasingly rely on machine learning models to make data-driven decisions based on interactions with users, infrastructure and other systems. For broader adoption, this practice must (*i*) accommodate product engineers without ML backgrounds, (*ii*) support finegrain product-metric evaluation and (*iii*) optimize for product goals. To address shortcomings of prior platforms, we introduce general principles for and the architecture of an ML platform, *Looper*, with simple APIs for decision-making and feedback collection. Looper covers the end-to-end ML lifecycle from collecting training data and model training to deployment and inference, and extends support to personalization, causal evaluation with heterogenous treatment effects, and Bayesian tuning for product goals. During the 2021 production deployment, Looper simultaneously hosted 440-1,000 ML models that made 4-6 million real-time decisions per second. We sum up experiences of platform adopters and describe their learning curve.

## 1 INTRODUCTION

With growing adoption of machine learning (ML), personalization is proving essential to competitive user experience [18]. To support users with different preferences, one needs good default tactics, user feedback, prioritizing delivered content and available actions [37]. When managing limited resources, e.g., for video serving, similar logic applies to network bandwidth, response latency, and video quality [23, 34]. This paper explores the use of ML for personalized decision-making in software products using what we call *smart strategies* (Section 2). Making smart strategies available to product engineers is challenging: (1) Long-term product objectives rarely match closed-form ML loss functions. (2) Product-generated data drifts away from training data. (3) Capturing correlations in training data (via ML) does not imply causal improvement of product metrics [10, 53]. (4) Clean data to evaluate the performance of ML systems is often unavailable, necessitating A/B testing. (5) Traditional A/B tests neglect personalized treatments. (6) Platforms to train, host and monitor hundreds of ML models are needed and promise economies of scale. (7) Real-time feature extraction and inference are needed, despite more efficient async batch processing. (8) Product engineers, many new to ML, need a simple, standard, future-proof way to embed smart strategies into products.

**Data-centric ML development** is a recent concept of refocusing ML development from models to data [35]. It supports software personalization with off-the-shelf models, where collecting the right data and selecting the appropriate class of models become primary differentiators [37]. Compared to developing and training ML models, data adequacy is often overlooked [45], and product platforms must use automation to compensate. Per Andrew Ng, "everyone jokes that ML is 80% data preparation, but no one seems to care" [44]. Yet, directly handling data sets and ML models in product code is cumbersome. Instead, *software-centric* ML integration with data collection and decision-making APIs offers a front-end to MLOps automation (Sections 2.2 and 3.2). Additionally, ML development often neglects structure in product evaluation data (Section 3.4).

**Vertical ML platforms** lower barriers to entry and support the entire lifecycle of ML models (Figure 1) in a repeatable way. Horizontal ML platforms provide storage, support data pipelines and offer basic services, whereas vertical platforms foster the reuse of not only ML components, but also workflows. At firms like Google, Meta, LinkedIn, Netflix, *specialized end-to-end vertical platforms* drive flagship product functionalities, such as recommendations. They have also been applied to software development, code quality checks, and even to optimize algorithms such as sorting and searching [15]. Platforms are built on ML frameworks like TensorFlow [1] and PyTorch [33] that focus on modeling for generic ML tasks, support hardware accelerators, and act as toolboxes for application development [24, 36]. Supporting smart strategies requires *general-purpose vertical platforms* to offer end-to-end ML lifecycle management. General-purpose vertical ML platforms can be internal to a company — Apple's Overton [41] and Uber's Michelangelo [27], — or broadly available to cloud customers — Google's Vertex, Microsoft's Azure Personalizer [2] and Amazon Personalize. A common theme is to help engineers "build and deploy deep-learning applications without writing code" via high-level, declarative abstractions [37]. Improving user experience and system performance with ML remains challenging [40] as correlations in data found by ML models might not lead to causal improvements. Little is known about optimizing for product goals [37, 52].

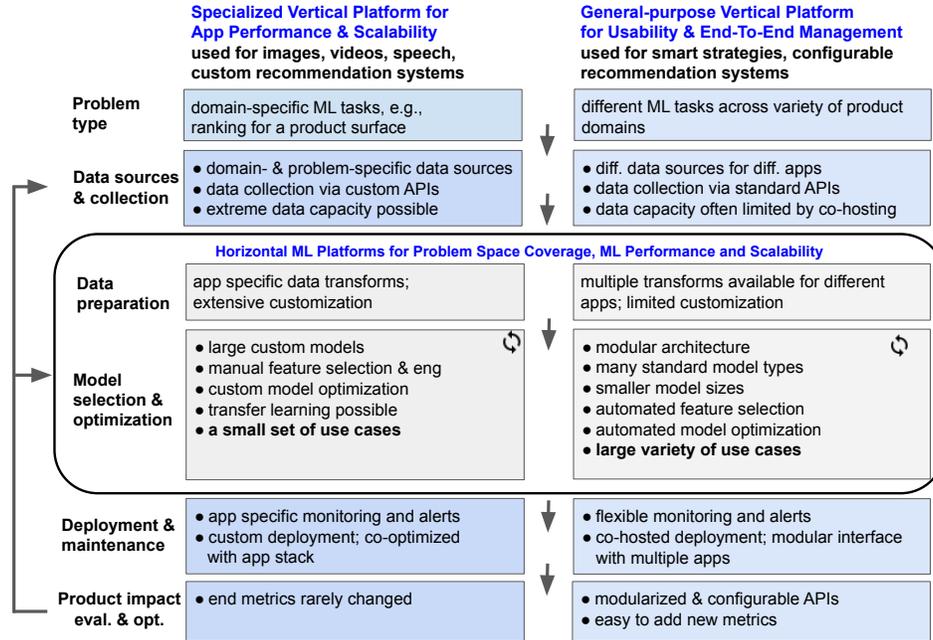| | **Specialized Vertical Platform for App Performance & Scalability** used for images, videos, speech, custom recommendation systems | **General-purpose Vertical Platform for Usability & End-To-End Management** used for smart strategies, configurable recommendation systems |
|---|---|---|
| **Problem type** | domain-specific ML tasks, e.g., ranking for a product surface | different ML tasks across variety of product domains |
| **Data sources & collection** | ● domain- & problem-specific data sources<br>● data collection via custom APIs<br>● extreme data capacity possible | ● diff. data sources for diff. apps<br>● data collection via standard APIs<br>● data capacity often limited by co-hosting |
| | **Horizontal ML Platforms for Problem Space Coverage, ML Performance and Scalability** | |
| **Data preparation** | app specific data transforms; extensive customization | multiple transforms available for different apps; limited customization |
| **Model selection & optimization** | ● large custom models<br>● manual feature selection & eng<br>● custom model optimization<br>● transfer learning possible<br>● **a small set of use cases** | ● modular architecture<br>● many standard model types<br>● smaller model sizes<br>● automated feature selection<br>● automated model optimization<br>● **large variety of use cases** |
| **Deployment & maintenance** | ● app specific monitoring and alerts<br>● custom deployment; co-optimized with app stack | ● flexible monitoring and alerts<br>● co-hosted deployment; modular interface with multiple apps |
| **Product impact eval. & opt.** | ● end metrics rarely changed | ● modularized & configurable APIs<br>● easy to add new metrics |

Fig. 1. Categories of applied ML platforms: horizontal vs. vertical, specialized vs. general-purpose (back arrows show vertical optimizations based on product metrics, see Section 3.3). Specialized platforms are limited in their support for diverse applications.

We develop support for *data-driven real-time smart strategies* via a *general-purpose vertical end-to-end* ML platform called *Looper*, internal to Meta, for rapid, low-effort deployment of moderate-sized models. Looper is a *declarative ML system* [27, 36, 37, 41] with coding-free full-lifecycle management of smart strategies via a GUI.

**Our technical contributions** include ①	a full-stack real-time ML platform (Section 3) with causal product-impact evaluation and optimization and handling of heterogeneous treatment effects (Sections 3.4, 4.2) via an experiment optimization system and meta-learners, ② a generic framework for *targeting long-term outcomes* by parameterized policies using plug-in supervised learning models and Bayesian optimization (Sections 3.3, 3.4, 4.2), ③ the *strategy blueprint* abstraction to optimize not only models, but the entire ML stack (Figures 1, 3 and Section 3.3), ④ capturing decision inputs and observations online via the succinct Looper API for product code; not only predicting what's logged by the API, but also optimizing black-box product objectives (Section 3.2), ⑤ broad deployment and substantial impact on product metrics (Section 4), ⑥ analysis of resource-usage bottlenecks (Appendix A), ⑦ qualitative analysis of our platform via a survey of adopters (Appendix B).

Specialized vertical ML platforms limit application diversity, while Looper hosts hundreds of production use cases thanks to its *general-purpose architecture*. Many vertical platforms [2, 15, 36, 41], don't solve as wide a selection of ML tasks as Looper does (classification, estimation, value and sequence prediction, ranking, planning) using supervised and reinforcement learning. Unlike platforms with asynchronous batch-mode feature extraction and inference [25, 27], Looper runs in real time and optimizes resource usage accordingly (Appendix A). To balance model quality, size and inference time,

Looper AutoML selects models and hyperparams, and performs vertical optimizations via *strategy blueprints* (Section 3.3).

**In the remainder of the paper**, Section 2 explores ML-driven smart strategies and relevant platform needs. Section 3 covers design principles for the Looper platform, introduces the architecture, the API, the blueprints, and specializations. Section 4 summarizes product impact at Meta, comparisons to baselines and adoption statistics. Appendices provide data to foster reproduciibility.

## 2 ML FOR SMART STRATEGIES

Compared to benchmark-driven research, ML that interacts with the world runs into additional challenges. In this paper, we target smart strategies at key decision points in software products, e.g.,

- application settings and preferences: selecting between defaults and user-specified preferences
- adaptive interfaces — certain options are shown only to users who are likely to pursue them
- controlling the frequency of ads, user notifications, etc
- prefetching or precomputation to reduce latency
- content ranking and prioritizing available actions

User preferences and context complicate decision-making. Simplifying a UI menu can boost product success, but menu preferences vary among users. Prefetching content to a mobile device can enhance user experience, but may require predicting user behavior.

While human-crafted heuristic strategies often suffice as an initial solution, ML-based smart strategies tend to outperform heuristics upon sufficient engineering investment [15, 29]. The Looper platform aims to lower this crossover point to broaden the adoption of smart strategies and deliver product impact over diverse applications. In this section, we discuss modeling approaches to enable smart strategies and cover the priorities in building such a platform.

## 2.1 Modeling approaches for smart strategies

Smart strategies are backed by supervised learning, contextual bandits (CB), and/or MDP-style reinforcement learning (RL). Given a model type, product decision problems need (1) ML optimization objectives to approximate the product goal(s) and (2) a decision policy to convert objective predictions into a single decision.

**Approximating product goals with predictable outcomes** (alternatively referred to as *proxy* or *surrogate* objectives) is a major difference between industry practice and research driven by existing ML models with abstract optimization objectives [48]. Good proxy objectives should be readily measurable and reasonably predictable. In recommendation systems, the "surrogate learning problem has an outsized importance on performance in A/B testing but is difficult to measure with offline experiments" [16]. We note a delicate tradeoff between easy-to-measure objectives directly linked to the decision vs. more complex objectives, e.g., ad clicks vs. conversions. Furthermore, product goals often implicitly have different weighting functions than the ML objective (e.g., the feedback provided by most prolific product users does not always represent other users [12]). Objectives can be modeled directly using *supervised learning*; alternatively, using CBs can model uncertainty in predictions across one or more objectives, which may then be used for exploring the set of optimal actions, e.g., in Thompson sampling [2, 3, 19, 32]. The use of RL enables the optimization of long-term, cumulative objectives, helping use cases with sequential dependencies [7, 24, 32]. To evaluate any one of these types of models and decision rules, true effects of the ML-based smart strategies can be estimated via A/B tests.

**Decision policies** postprocess the raw model outputs into a final product decision or action. For single-objective tasks in supervised learning this may be as simple as making a binary decision if the objective prediction exceeds a threshold, e.g. turning the probability of a click into a binary prefetch decision (Section 4.3). For tasks with multiple objectives and more complex action spaces, the template for a decision policy is to assign a scalar value or score to all possible *actions* in the decision space, which can then be ranked through sorting. In recommendation systems, a standard approach is to use a combination function (usually a weighted product of objective predictions) to generate a score for each candidate [55]. When using reinforcement learning, reward shaping [31] weighs task scores in the reward function to optimize for the true long-term objective. Optimizing this weighting for multi-objective tasks is explored in Section 3.3. More sophisticated policies also use randomization to explore the action space, e.g. Thompson sampling in contextual bandits [19], or $\varepsilon$-greedy approaches for exploration in ranking [3].

## 2.2 Extending end-to-end ML for smart strategies

Traditional end-to-end ML systems go as far as to cover model publishing and serving [27, 36, 37, 41], but to our knowledge rarely track *how* the model is used in the software stack. Assessing and optimizing the impact of smart strategies, especially with respect to product goals, requires experimentation on all aspects of the modeling framework – from metric and model selection to policy optimization. To streamline this experimentation and reap its benefits, smart-strategies platforms must extend the common definition of end-to-end into the software layer.

**Software-centric ML integration** [2, 15] – where data collection and decision-making are fully managed through platform APIs – enables both high-quality data collection and holistic experimentation. Notably, the platform can now keep track of all decision points and support A/B tests between different configurations. Well-defined APIs improve adoption among product engineers with limited ML background, and ML configuration can be abstracted via declarative programming or GUI without requiring coding [37].

**End-to-end AutoML.** Hyperparameter tuning is often automated per model via black-box optimization [11]. But optimizing the loss function of SOTA models by 1% often brings no long-term product gains, whereas tuning decision policy params usually helps, e.g., by better reflecting penalties for false positives/negatives. In our *full-stack* (extended end-to-end) regime, we enable AutoML for the entire ML pipeline via declarative *strategy blueprints* (Section 3.3) and an adaptive experiments framework tied to product metrics [9].

## 2.3 Additional requirements for smart strategies

**Metadata features** for product-specific models (e.g., account type, time spent online, interactions with other accounts) introduce new aspects to learning smart strategies in addition to traditional content features (images, text, video) commonly handled by ML platforms. Unlike image pixels, metadata features are diverse, uncorrelated, require non-uniform preprocessing, and are often *joined* from different sources. Patterns in metadata change quickly, necessitating regular retraining of ML models on fresh data, as well as monitoring and alerts. Interactions between dense metadata features can often be handled by GBDTs or shallow neural nets. Sparse and categorical features need adequate representations [42] and special provisions if used by neural network architectures [38].

**Non-stationary environments** are typical for deployed products but not for research prototypes and SOTA results.

**Logging and performance monitoring** are important capabilities for a production system. Dashboards monitor system health and help understand model performance in terms of statistics, distributions and trends of features and predictions, automatically triggering alerts for anomalies [4, 13]. Our platform integrates with Meta 's online experimentation framework, and production models can be withdrawn quickly if needed.

**Monitoring and optimizing resource usage** flags inefficiences across training and inference. Our monitoring tools track resource usage to components of the training and inference pipeline (Section 3.2), and help trade ML performance for resources and latency.

## 3 THE LOOPER PLATFORM

To support a smart strategy, a vertical ML platform (Figure 1) collects features and labels from a running product, trains a model, and produces predictions in real time for use in the product. Such "loops" need operational structure — established processes and protocols for model revision and deployment, evaluation and tracking of product impact, and overall maintenance. We now introduce insights, design principles and an architecture for a vertical smart-strategies platform to address the needs outlined in Sections 1 and 2.

## 3.1 Design principles and a concept inventory

In contrast to heavy-weight ML models for vision, speech and NLP that favor offline inference (with batch processing) and motivate applications built around them, we address the demand for smart strategies within software applications and products. These smart strategies operate on metadata — a mix of categorical, sparse, and dense features, often at different scales. Respective ML models are lightweight, they can be re-trained regularly and deployed quickly on shared infrastructure in large numbers. Downside risks are reduced via (*i*) simpler data stewardship, (*ii*) tracking product impact, (*iii*) failsafe mechanisms to withdraw poorly performing models. Smart strategies have a good operational safety record and easily improve naive default behaviors.

The human labeling process common for CV and NLP fails for metadata because relevant decisions and predictions (a) only make sense in an application context, (b) in cases like *data prefetch* (Section 4.3) only make sense to engineers, (c) may change seasonally, and even daily. Instead of human labeling, our platform interprets user-interaction and system-interaction metadata as either labels for supervised learning or rewards for reinforcement learning. To improve operatonal safety and training efficiency, we rely on batch-mode (offline) training, even for reinforcement learning. Given real-time inference, model agility beyond daily re-training is supported by real-time engineered features, such as event counters.

Our platform ensures fast onboarding, robust deployment and low-effort maintenance of multiple smart strategies where positive impacts are measured and optimized directly in application terms (Appendix B). To this end, we separate application code from platform code, and leverage existing horizontal ML platforms with interchangeable models for ML tasks (Figure 1). Intended for company engineers, our platform benefits from high-quality data and engineered features in the company-wide feature store [39]. To simplify onboarding for product teams and keep developers productive, we automate and support

- Workflows avoided by engineers [45], e.g., feature selection and preprocessing, and tuning ML models for metadata.
- Workflows that are difficult to reason about, e.g., tuning ML models to product metrics.

We first introduce several concepts for platform design.

The **decision space** captures the shape of decisions within an application which can be made by a smart strategy. It can be just {0,1} to indicate whether a notification is shown. It can be a continuous-value space for time-to-live (TTL) of a cache entry. It can be a data structure with configuration values for a SW system, such as a live-video stream encoder. With reinforcement learning, the decision space matches well with the concept of action space.

**Application context** captures necessary key information provided by a software system at inference time to make a choice in the decision space. The application context may be directly used as features or it may contain ID keys to extract the remaining features from the feature store (Section 3.3).

**Product metrics** evaluate the performance of an application and smart strategies. When specific decisions can be judged by product metrics, one can generate labels for supervised learning, unlike for metrics that track long-term objectives.

**A proxy ML task** casts product goals in mathematical terms to enable (*i*) reusable ML models that optimize formal objectives and (*ii*) decision rules that map ML predictions into decisions (Section 2.1). Setting proxy tasks draws on domain expertise, but our platofrm simplifies this process.

**Evaluation of effects on live data** verifies that solving the proxy task indeed improves product metrics. Access to Meta's monitoring infrastructure helps detect unforeseen side effects. As in medical trials, (1) we need evidence of a positive effect, (2) side-effects should be tolerable, and (3) we should not overlook evidence of side-effects. On our platform, product developers define the decision space, allowing the platform to automatically select model type and hyper-parameter settings. The models are trained and evaluated on live data without user impact, and improved until they can be deployed. Newly trained models are *canaried* (deployed on shadow traffic) before product use – such models are evaluated on a sampled subset of logged features and observations, and offline quality metrics (e.g., MSE for regression tasks) are computed. This helps avoid degrading model quality when deploying newer models.

## 3.2 Platform architecture: the core

Traditional ML pipelines build training data offline, but our platform uses a *live feature store* and differs in two ways:

- **Software-centric vs. data-centric interfaces.** Rather than passed via files or databases, training data are logged from product surfaces as Looper APIs intercept decision points in product software. Product engineers delegate training-data quality concerns (missing or delayed labels, etc) to the platform. Missing data are represented by special values.
- **An online-first approach.** Looper API logs live features and labels at the decision and feedback points, then joins and filters them via real-time stream processing. This *immediate materialization* avoids data hygiene issues [2] and storage overhead: it keeps training and inference consistent and limits label leakage by separating features and labels in time. Looper's *complete chain of custody for data* (without exposing data tables or files) helps prevent engineering mistakes.

**The Looper RPC API** relies on two core methods:

**I.** `getDecision(decision_id, application_context)` returns a decision-space value, e.g., `True`/`False` for binary choices or a floating-point score for ranking. Unlike in the 3-call APIs in [2, 15],
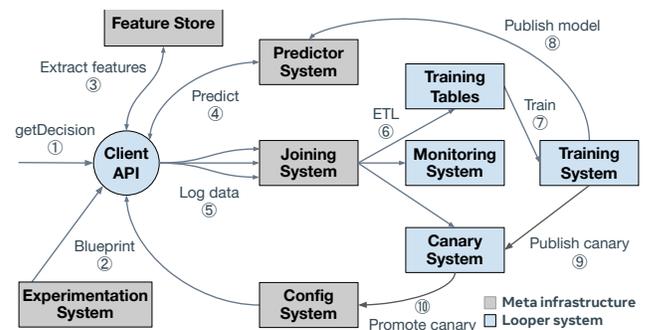


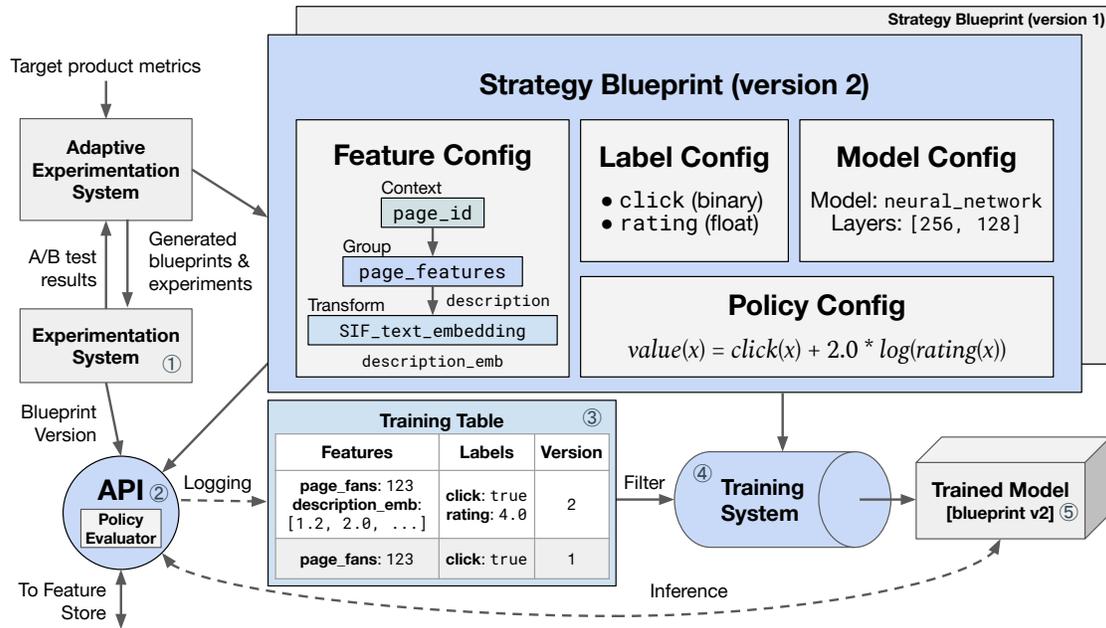Fig. 2. Data flow in the Looper platform. Figure 3 expands the left side.

Fig. 3. The strategy blueprint and how it controls different aspects of the end-to-end model lifecycle. Continuation of Figure 2.

null is returned before a model is available (to trigger default behavior). User-defined `decision_id` ties each decision with observation(s) logged later (II); it may be randomly generated for clients to use. `application_context` is a dictionary representation of the application context (Section 3.1), e.g., with the user ID (used to retrieve additional user features), current date/time, etc.

**II.** | `logObservations(decision_id, observations)` | logs labels for training proxy ML task(s), where `decision_id` must match a prior `getDecision` call. Observations capture users' interactions, responses to a decision (e.g., clicks or navigation actions), or environmental factors such as compute costs.

Though deceptively simple in product code, this design fully supports the MLOps needs of the platform. We separately walk through the online (inference) and offline (training) steps of the pipeline in Figure 2. ① Product code initializes the Looper client API with one of the known strategies registered in the UI. `getDecision()` is then called with the `decision_id` and `application_context`. ② Looper client API retrieves a versioned configuration (the "strategy blueprint", Section 3.3) for the strategy to determine the features, the model instance, etc. The exact version used may be controlled through an external experimentation system. ③ The client API passes the application context to the Meta feature store (Section 3.3), which returns a complete feature vector. ④ The client API passes the feature vector and production model ID to a distributed model predictor system (cf. [47]), which returns proxy task predictions to the client. Then, the client API uses a decision policy (Section 2.1) to make the final decision based on the proxy predictions. Decision policies are configured in a domain-specific language (DSL) using logic and formulas. ⑤ Asynchronously, the anonymized feature vector and predictions are logged to a distributed online joining system (c.f.

[5]), keyed by the decision ID and marked with a configurable and relatively short TTL (time-to-live). The `logObservations` API also sends (from multiple request contexts) logs to this system. Complete "rows" with matching features and observations are logged to a training table, with retention time set per data retention policies. The remaining steps are performed offline and asynchronously.

⑥ Delayed and long-term observations are logged in a table and then joined offline via Extract, Transform, and Load (ETL) pipelines [6]. These pipelines perform complex data operations such as creating MDP sequences for reinforcement learning. The logged features, predictions, and observations are sent for logging and real-time monitoring as per Section 2.2. ⑦ An offline training system [22] retrains new models nightly or when sufficient data are available, addressing concerns from Section 3.1. ⑧ Trained models are published to the distributed predictor for online inference. ⑨ Models are then registered for canarying (Section 3.1). ⑩ A canary model that outperforms the prior model is promoted to production and added to the loop configuration.

### 3.3 Product optimization with strategy blueprints

The end-to-end nature of the Looper platform brings its own set of challenges regarding data and configuration management in the system. Existing ML management solutions [49] primarily focus on managing or versioning of data and models, which is insufficient in covering the full lifecycle of smart strategies. In this section we introduce the concept of a *strategy blueprint*, a version-controlled configuration that describes how to construct and evaluate a smart strategy. Blueprints are immutable, and modifications (typically through a GUI) create new versions that can be compared in production through an online experimentation platform, allowing for easy

rollback if needed. The strategy blueprint (Figure 3) controls four aspects of the ML model lifecycle and captures their cross-product:

**Feature configuration.** Modern ML models can use thousands of features and computed variants, which motivates a unified repository — a *feature store* — usable in model training and real-time inference [26, 39]. Features within a *feature group* are associated with the same application context (e.g., a Web page id) and are computed together. Feature variants are produced by *feature transforms*, e.g., pre-trained or SIF [8] text embeddings. The Looper blueprint leverages feature stores for feature management and tracks (*i*) feature groups via a computational graph and (*ii*) downstream feature transforms. Blueprint modifications often try to improve model quality by experimenting with new features.

**Label configuration** controls how ML objectives are proxied by clicks, ratings, etc, as per Section 2.1. Faithful proxies for product metrics are hard to find [48], hence experimentation with label sets.

**Model configuration** helps product teams explore model architecture tradeoffs (DNNs, GBDTs, reinforcement learning). The blueprint only specifies high-level architecture parameters, while lower-level hyperparameters (e.g., learning rates) are delegated to AutoML techniques invoked by the training system (Section 2.2).

**Policy configuration** controls how raw objective predictions are translated into decisions (Section 2.1). It uses a lightweight domain specific language (DSL). Figure 3 illustrates a ranking decision, where the click and rating objectives are weighted and combined to generate a single score per candidate. Smart strategies often need to optimize the importance weights embedded in decision policies.

Blueprints capture *compatibility* between versions, e.g., the training pipeline for ver. *A* may use data from ver. *B* if features and labels in *A* are subsets of those in *B*. Tagging each training row with the originating blueprint version enables data sharing between versions.

Figure 3 illustrates the lifecycle of a blueprint. From left to right: ① An experimentation system enables different blueprint versions to be served across the user population to facilitate A/B testing (optionally, in concert with a "blueprint optimizer", described later below). ② The client API uses the blueprint feature configuration to obtain a complete feature vector from the feature store. ③ Completed training examples are logged to training tables, tagged with the originating blueprint version. ④ The training system filters data by compatible version and executes the pipeline per the blueprint's feature, label, and model configurations. The policy configuration may be needed as well for more sophisticated model types (reinforcement learning). ⑤ Trained models are published in the blueprint version. So, the client API uses only models explicitly linked to its served blueprint version. To generate the final product-facing decision, the client also uses the policy configuration.

**Vertical optimizations** with the blueprint abstraction capture dependencies between the four configuration types. Since long-term product metrics are rarely available in closed form, such an optimization requires (*i*) a sequence of A/B tests that evaluate product metrics and (*ii*) parameter adjustment between these tests. Given that each A/B test (experiment) can take significant time and impact many end-users, very few A/B tests can be used in practice. This calls for (multi-objective) Bayesian optimization: to tune parameters in a blueprint and optimize product metrics w/o closed-form

representation, Looper leverages an *adaptive experimentation* platform [9], see Section 4.2 for details. Product outcomes often improve even by just tuning weights in the policy configuration, e.g., for recommendation scores and reward shaping.

**Bookkeeping** for product groups using Looper is performed in terms of *use cases* tied to an *application context* and an *ML task* (Section 3.1). Multiple candidate configurations are maintained to train and evaluate candidate *model instances* and define *decision policies*. While only one blueprint (and model) per use case is in production, many model instances may be undergoing shadow evaluation.

## 3.4 Platform specializations

The core platform (Sections 3.2 and 3.3) goes a long way to address the challenges listed in the Introduction. However, additional structures are needed because (*i*) product-metric evaluation and optimization have serious blind spots, while (*ii*) some classes of application are cumbersome to support. Platform specializations that address these deficiencies add significant value to the platform.

**Integrated experiment optimizations.** Even when a product metric is approximated well by an ML loss function, the correlations captured by the model might not lead to causal product improvements. Hence, A/B testing estimates the *average treatment effect* (ATE) of the change across product users. Shared repositories of product metrics are common [10, 28, 53], and product variants are systematically explored by running many concurrent experiments [9]. While dealing with non-stationary measurements, balancing competing objectives, and supporting the design of sequential experiments [9], a common challenge with A/B tests is to find subpopulations where treatment effects differ from the global ATE – *heterogeneous treatment effects* (HTE). Common neglect for HTEs in A/B testing leaves room for improvement [10, 12] [50], likely delivering suboptimal treatments. The Looper platform and its support for A/B testing dramatically simplify HTE modeling on the Meta online experimentation platform, and help deploying treatment assignments based on HTE estimates.

In an initial training phase, Looper's `getDecision()` API acts as a drop-in replacement for the standard A/B testing API, and falls through to a standard randomized assignment while still logging features for each experiment participant. Then, metrics from the standard A/B testing repertoire help derive the treatment outcome (observations) for each participant, and the Looper platform trains specialized HTE models (meta-learners such as T-, X-, and S- learners [30]). In a final step, the HTE model predictions can be used in a decision policy to help `getDecision()` make intelligent treatment assignments and measurably improve outcomes compared to any individual treatment alone. In this scenario, the best HTE estimate for a given user selects the actual treatment group. Our integration links Looper to an established experiment optimization system [9] and creates synergies discussed in Section 4.2. A further extension relaxes the standard A/B testing contract to support *fully dynamic assignments* and enables reinforcement learning [7].

**Looper for ranking.** The `getDecision()` + `logObservations()` API is general enough to implement simple recommendation systems, but advanced systems need finer support. Higher-ranked items
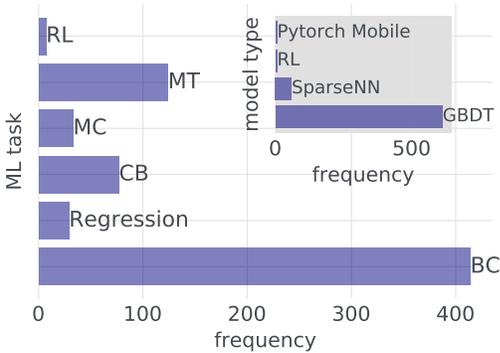
Fig. 4. ML tasks and model types on our platform. On the outer plot — ML tasks: Reinforcement Learning (RL), multitask (MT), multiclass (MC) and binary classification (BC), contextual bandit (CB), regression. Multitask models blend classification or regression sub-tasks. On the inner plot — model types: PyTorch Mobile, RL, NN, Gradient Boosted Decision Trees.

are more often chosen by users, and this *positional bias* can be handled (in the API) by including the displayed position as a special input during training [17]. To derive a final priority score for each item, the multiple proxy task predictions are often combined through a weighted combination function [55]. Recommender systems learn from user feedback, as long ass lesser-explored items are occassionally included among top results (the explore/exploit tradeoff [54]). A specialized Looper ranking system abstracts these considerations in a higher-level API (getRanking) allowing the ordering of an entire list of application contexts, and also allows recording of display-time observations such as the relative screen position of each item.

## 4 PRODUCTION DEPLOYMENT

Looper supports real-time inference with moderate-sized models to improve various aspects of software systems. These models are trained and deployed quickly and maintained on the platform, whereas our two-call RPC API (Section 3.2) decouples platform code from application code. Deployed in production at Meta during the entire year 2021, Looper improved product metrics, made use of compute resources more efficient, and streamlined maintenance. To enhance the reproducibility of our work, we describe empirical observations and prominent applications, outline adoption and impact, and summarize adopter survey results (in Appendix B). Resource estimates are in "server" units (we do not use GPUs). For smart strategies that learn from end-users, performance and impact inevitably depend on the user base and product metrics infrastructure. Hence, we report Looper impact as a fraction of product teams's half-year results. Resource savings are reported as percent vs. baselines.

### 4.1 Statistics for ML tasks and models in Looper

Given nonstationary application environments, Looper retrains models on the day when 20% fresh training data become available (or on a set schedule). Choosing appropriate ML models requires trading off performance metrics with resource usage, inference latency, and configuration effort. SVM packages struggle with multimodal data and scale poorly to voluminous data. DNNs scale to 1B+ data rows, handle sparse features, but use more memory than simpler models. DNNs are sensitive to architecture configuration and are

currently less explored for tabular data, making configuration challenging without ML expertise. Gradient-Boosted Decision Trees (GBDT, XGBoost) are compact and robust, handle multimodal tabular data naturally (including sparse features and missing values), do not require architecture search or GPUs, and scale to 100M rows. Typical *inference latency* with GBDT and XGBoost is in low single *ms* for server-side models and 1-2 ms for (smaller) mobile models. Inference with DNN-based models, especially those using latent embeddings, is an order of magnitude slower. That's why the mean inference latency (10ms) across Looper use cases is greater than the median (2ms). 90%le and 99%le latencies are 20ms. Figure 4 summarizes ML tasks deployed on Looper and the models selected for them. Figure 6 shows that models typically use 50-200 features, and most features are used by many models. Feature extraction latency has median 45ms and mean 120ms. Latencies for extr†acting synthetic/engineered features are greater (90%le and 99%le is 240 ms). For a broader picture, Figure 5 summarizes 29 utilization rates for system resources. In particular, feature extraction tends to be a bottleneck. This data is further discussed in Appendix A and used to optimize resource utilization in Looper.

### 4.2 Application deep dive – personalized experiments

Section 3.4, while focusing on our platform architecture, outlined a special application of smart strategies — embedding them into the standard experimentation framework to automate the personalization of A/B treatments when HTEs are detected and captured by ML models. As A/B testing APIs are common and accessible [10, 53], Looper's integration with the standard experimentation APIs makes training and deploying personalized smart strategies as easy as changing *one or two lines* of code by product teams. In practice, exposing smart strategies through such APIs brings several benefits:

- Client code and the learning curve are simplified by repurposing the decision API as the A/B testing API.
- Dataset preparation and modeling flow are automated for the task of optimizing metric responses based on end-users exposed to each treatment. Metric responses can be automatically sourced from the experimentation measurement framework without manual labeling.
- Product metrics are traded off across many strategies, offline and online, via multi-objective optimization (MOO) [9, 20].
- Smart strategies are automatically compared against all baseline treatments (i.e., A or B), making the tradeoff between metric impact and costs explicit.

Previously such experiment optimization needed dedicated engineering resources. Now the tight integration of the Looper platform with the experimentation framework allows product engineers quickly evaluate a smart strategy and optimize its product impact *in several weeks*. With automatic MOO, engineers find tradeoffs appropriate to a given product context. For example, during a server capacity crunch, one team traded a slight deterioration in a product metric for a 50% resource savings. Predicating product deployment on such experiments creates safeguards against ML models that generalize poorly to live data. This also helps tracking product impact. For example, a user authentication use case [7] reduced SMS cost by 5% while remaining neutral for engagement metrics.

## 4.3 Application deep dive – data prefetching

Online applications strive to reduce response latency for user interactions. Optimized resource prefetching based on user history helps by proactively loading application data. Modern ML methods can accurately predict the likelihood of data usage, minimizing unused prefetches. Our Looper platform supports prefetching strategies for many systems within Meta, often deeply integrated into the product infrastructure stack. For example, Meta's GraphQL [14] data fetching subsystem uses our platform to decide which prefetch requests to service, saving both client bandwidth and server-side resources. It yields around 2% compute savings at peak server load. As another example, Meta's application client for lower-end devices (with a "thin-client" server-side rendering architecture [43]) also uses our platform to predictively render entire application screens. Our automated end-to-end system helps deploying both models and threshold-based decision policies then tune them for individual GraphQL queries or application screens, with minimal engineering effort. Based on numerous deployed prefetch models, we have also developed large-scale modeling of prefetching. User-history models have proven to be helpful for this task [51]; building up on this idea, we created application-independent vector embeddings based on users' surface-level activity. To accomplish this, we train a multi-task auto-regressive neural-network model to predict how long a user will stay in each application surface (e.g., news feed, search, notifications), based on a sequence of (`application surface, duration`) events from the user's history. As is common in CV and NLP, intermediate-layer outputs of this DNN predict prefetch accesses well and make specialized features unnecessary. Optimized prefetching illustrates how secondary, domain-specific platforms are enabled by the core Looper platform; infrastructure teams only need to wire up the prediction and labeling integration points while Looper provides full ML support.

## 4.4 Adoption and impact

Several internal vertical platforms at Meta [26] compete for a rich and diverse set of applications. Product teams sometimes relocate their ML models to a platform with greater advantages, while a few high-value applications are run by dedicated infrastructure teams. Looper was chosen and is currently used by 90+ product teams at Meta. On any day in 2021, these teams deployed 440-1K models that made 4-6 million decisions per second. Application use cases fall into five categories in decreasing order of usage (Figures 5 and 7):

- **Personalized Experience** is tailored based on the user's engagement history. For example, we display a new feature prominently only to those likely to use it.
- **Ranking** orders items to improve user utility, e.g., to personalize a feed of candidate items for the viewer.
- **Prefetching/precomputing** data/resources based on predicted likelihood of usage (Section 4.3).
- **Notifications/prompts** can be gated on a per-user basis, and sent only to users who find them helpful.
- **Value estimation** predicts regression tasks, e.g., latency or memory usage of a data query.

The impact of ML performance on product metrics varies by application. For a binary classifier, increasing ROC AUC from 90% to 95%
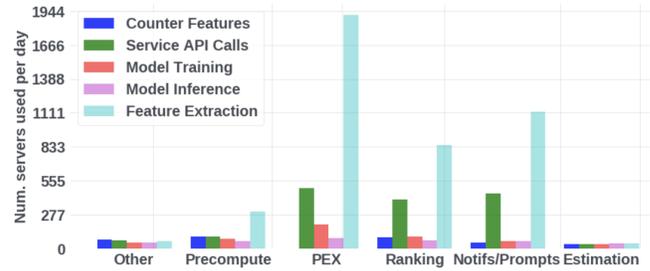


Fig. 5. Resource utilization rate by resource type and application category (see Section 4.4). The Service API category includes API calls other than feature extraction and prediction service.

might not yield large product gains when such decisions contribute little to product metrics if bottlenecks lie elsewhere. But increasing ROC AUC from 55% to 60% is impactful when each percent translates into tangible resource savings or other metrics, as it would be for online payment processing. Looper use cases contributed to compute savings (server utilization), user engagement (e.g., daily active users) and other top-line company reporting metrics. Many product teams at Facebook and Instagram adopted Looper without additional staffing, and it is common for Looper to contribute 20-40% of improvements to product goal metrics. In several cases, Looper helped product teams outperform their goals by over 2x.

## 5 CONCLUSIONS

We outline opportunities to embed data-driven self-optimizing smart strategies for product decisions into software systems, so as to enhance user experience, optimize resource utilization, and support new functionalities. We describe the deployment of smart strategies (at an unprecedented scale) through software-centric ML integration where decision points are intercepted and data is collected through APIs [2]. This process requires infrastructure and automation to reduce operational mistakes and maintain ML development velocity.

Our ML platform Looper addresses the complexities of product-driven end-to-end ML systems and facilitates at-scale deployment of smart strategies through technical insights, platform-level abstractions, a novel architecture, the use of Bayesian optimization, and interfaces with an adaptive experimentation system. As an important simplification, inference input processing matches that for training. The Looper product RPC API is simplified down to two calls. The Looper platform treats end-to-end ML development more broadly than prior work [37, 52], providing extensive support for product impact evaluation of smart strategies via causal inference. Looper learns heterogenous treatment effects (HTE) from product evaluation data by repurposing the Looper RPC API as a drop-in replacement for a standard A/B testing API. Vertical optimizations with long-term product objectives are enabled by our *strategy blueprint abstraction* and the use of Bayesian optimization.

As observed during production deployment in 2021, Looper offers immediate, tangible benefits in terms of data availability, easy configuration, judicious use of available resources, reduced engineering effort, and ensuring product impact. It makes smart strategies easily accessible to product engineers at large scale and enables product teams to build, deploy and improve ML-driven capabilities in a self-serve fashion without ML expertise. We observed product teams

launch smart strategies within their products in one month (Appendix B). The lower barriers to entry and faster deployment lead to more pervasive use of ML to optimize user experience in new products and old products not designed with ML in mind. Support for prefetching and personalized A/B testing have been in demand, whereas end-to-end management enables holistic resource accounting and optimization [52]. The overall impact in product metrics and resource-efficiency is substantial. We also provide empirical insights into usage by resource types and application category.

Long-term benefits of our platform approach include effort and module reuse, end-to-end reproducibility, consistent reporting, reliable maintenance, and being able to upgrade ML libraries and offer new ML model types with consistent interface. Successful Looper adopters often launch additional data-driven smart strategies, and this virtuous cycle encourages designing SW systems with built-in ML to enhance user experience and adaptation to the environment.

## REFERENCES

[1] Martín Abadi et al. 2016. TensorFlow: A System for Large-Scale Machine Learning. In *OSDI*. 265–283.

[2] Alekh Agarwal et al. 2016. Making contextual decisions with low technical debt. *arXiv:1606.03966* (2016).

[3] Deepak Agarwal, Bee-Chung Chen, and Pradheep Elango. 2009. Explore/exploit schemes for web content optimization. In *ICDM 2009*. 1–10.

[4] Saleema Amershi et al. 2019. Software engineering for machine learning: A case study. In *ICSE-SEIP*. 291–300.

[5] Rajagopal Ananthanarayanan et al. 2013. Photon: Fault-Tolerant and Scalable Joining of Continuous Data Streams. In *ACM SIGMOD Int'l Conf. on Management of Data (SIGMOD '13)*. ACM, 577–588. https://doi.org/10.1145/2463676.2465272

[6] Anonymous. 2021. ETL vs ELT: Must Know Differences. https://www.guru99.com/etl-vs-elt.html

[7] Pavlos Athanasios Apostolopoulos et al. 2021. Personalization for Web-based Services using Offline Reinforcement Learning. *arXiv:2102.05612* (2021).

[8] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2017. A Simple but Tough-to-Beat Baseline for Sentence Embeddings. *ICLR* (2017).

[9] Eytan Bakshy et al. 2018. AE: A domain-agnostic platform for adaptive experimentation. *NeurIPS 2018 Systems for ML Workshop*.

[10] Eytan Bakshy, Dean Eckles, and Michael S Bernstein. 2014. Designing and deploying online field experiments. In *WWW' 14*. 283–292.

[11] Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. 2020. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *NeurIPS 33*.

[12] Alex Beutel, Ed H Chi, Zhiyuan Cheng, Hubert Pham, and John Anderson. 2017. Beyond globally optimal: Focused learning for improved recommendations. In *WWW' 17*. 203–212.

[13] Eric Breck, Shanqing Cai, Eric Nielsen, Michael Salib, and D. Sculley. 2017. The ML Test Score: A Rubric for ML Production Readiness and Technical Debt Reduction. In *IEEE Big Data*.

[14] Lee Byron. 2015. GraphQL: A data query language. https://engineering.fb.com/2015/09/14/core-data/graphql-a-data-query-language

[15] Victor Carbune, Thierry Coppey, Alexander Daryin, Thomas Deselaers, Nikhil Sarda, and Jay Yagnik. 2018. SmartChoices: Hybridizing programming and machine learning. *arXiv:1810.00619* (2018).

[16] P. Covington, J. Adams, and E. Sargin. 2016. Deep neural networks for Youtube recommendations. In *RecSys*.

[17] Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. An experimental comparison of click position-bias models. In *WSDM 2008*. 87–94.

[18] Marc D'Arcy. 2021. Opinion: The 3 Post-COVID Trends Empowering People and Shaping the Future. https://adage.com/article/opinion/opinion-3-post-covid-trends-empowering-people-and-shaping-future/2342861

[19] Samuel Daulton et al. 2019. Thompson sampling for contextual bandit problems with auxiliary safety constraints. *arXiv:1911.00638* (2019).

[20] Samuel Daulton, Maximilian Balandat, and Eytan Bakshy. 2021. Parallel Bayesian Optimization of Multiple Noisy Objectives with Expected Hypervolume Improvement, In NeurIPS 34. *arXiv:2006.05078*.

[21] Ben Dickson. 2021. Why machine learning strategies fail. https://venturebeat.com/2021/02/25/why-machine-learning-strategies-fail/

[22] Jeffrey Dunn. 2016. Introducing FBLearner Flow: Facebook's AI backbone. https://engineering.fb.com/2016/05/09/core-data/introducing-fblearner-flow-facebook-s-ai-backbone

[23] Qing Feng, Benjamin Letham, Hongzi Mao, and Eytan Bakshy. 2020. High-dimensional contextual policy search with unknown context rewards using Bayesian optimization. *NeurIPS 33* (2020).

[24] Jason Gauci et al. 2018. Horizon: Facebook's Open Source Applied Reinforcement Learning Platform. *arXiv:1811.00260* (2018).

[25] Udit Gupta et al. 2020. The Architectural Implications of Facebook's DNN-based Personalized Recommendation. *HPCA* (2020), 488–501. arXiv:1906.03109

[26] Kim Hazelwood et al. 2018. Applied machine learning at Facebook: A datacenter infrastructure perspective. In *HPCA 2018*. IEEE, 620–629.

[27] Jeremy Hermann and Mike Del Balso. 2017. Meet Michelangelo: Uber's Machine Learning Platform. https://eng.uber.com/michelangelo-machine-learning-platform/

[28] Ron Kohavi, Roger Longbotham, Dan Sommerfield, and Randal M Henne. 2009. Controlled experiments on the Web: survey and practical guide. *Data mining and knowledge discovery* 18, 1 (2009), 140–181.

[29] Tim Kraska et al. 2017. The Case for Learned Index Structures. *CoRR* (2017). arXiv:1712.01208

[30] Sören R. Künzel et al. 2019. Metalearners for estimating heterogeneous treatment effects using machine learning. *PNAS* 116, 10 (Feb 2019), 4156–4165.

[31] Adam Daniel Laud. 2004. *Theory and application of reward shaping in reinforcement learning*. UIUC.

[32] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. 2010. A contextual-bandit approach to personalized news article recommendation. In *WWW*. 661–670.

[33] Shen Li et al. 2020. PyTorch Distributed: Experiences on Accelerating Data Parallel Training. In *VLDB*, Vol. 13(12).

[34] Hongzi Mao et al. 2020. Real-world video adaptation with reinforcement learning. *arXiv:2008.12858* (2020).

[35] Lester James Miranda. 2021. Towards data-centric machine learning: a short review. (2021). https://ljvmiranda921.github.io/notebook/2021/07/30/data-centric-ml/

[36] P. Molino, Y. Dudin, and S. S. Miryala. 2019. Ludwig: a type-based declarative deep learning toolbox. *arxiv:1909.07930* (2019).

[37] P. Molino and C. Ré. 2021. Declarative Machine Learning Systems. *ACM Queue* 19 (2021). Issue 3.

[38] Maxim Naumov et al. 2019. Deep Learning Recommendation Model for Personalization and Recommendation Systems. *CoRR* abs/1906.00091 (2019).

[39] Laurel J. Orr et al. 2021. Managing ML Pipelines: Feature Stores and the Coming Wave of Embedding Ecosystems. *CoRR* (2021). arXiv:2108.05053

[40] Andrei Paleyes, Raoul-Gabriel Urma, and Neil D Lawrence. 2020. Challenges in deploying machine learning: a survey of case studies. *arXiv:2011.09926* (2020).

[41] Christopher Ré et al. 2019. Overton: A data system for monitoring and improving machine-learned products. *arXiv:1909.05372* (2019).

[42] Pau Rodríguez, Miguel A Bautista, Jordi Gonzàlez, and Sergio Escalera. 2018. Beyond One-hot Encoding: lower dimensional target embedding. *arXiv:1806.10805* (2018).

[43] Gautam Roy. 2016. How we built Facebook Lite for every Android phone and network. https://engineering.fb.com/2016/03/09/android/how-we-built-facebook-lite-for-every-android-phone-and-network

[44] Ram Sagar. 2021. Andrew Ng Urges ML Community To Be More Data-Centric. https://analyticsindiamag.com/big-data-to-good-data-andrew-ng-urges-ml-community-to-be-more-data-centric-and-less-model-centric/

[45] Nithya Sambasivan et al. 2021. "Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes AI. *SIGCHI, ACM* (2021).

[46] David Sculley et al. 2015. Hidden technical debt in machine learning systems. *NIPS* 28 (2015), 2503–2511.

[47] Jonathan Soifer et al. 2019. Deep Learning Inference Service at Microsoft. In *USENIX Conf. Operational Machine Learning (OpML)*. USENIX, 15–17. https://www.usenix.org/conference/opml19/presentation/soifer

[48] Gregory J. Stein. 2019. Proxy metrics are everywhere in machine learning. http://cachestocaches.com/2019/1/proxy-metrics-are-everywhere-machine-lea

[49] M. Vartak and S. Madden. 2018. MODELDB: Opportunities and Challenges in Managing Machine Learning Models. *IEEE Data Eng. Bull.* 41, 4 (2018), 16–25.

[50] S. Wager and S. Athey. 2018. Estimation and inference of heterogeneous treatment effects using random forests. *J. Amer. Stat. Assoc.* 113, 523 (2018), 1228–1242.

[51] Hanson Wang, Zehui Wang, and Yuanyuan Ma. 2019. Predictive Precompute with Recurrent Neural Networks. *arXiv:1912.06779* (2019).

[52] Carole-Jean Wu et al. 2021. Sustainable AI: Environmental Implications, Challenges and Opportunities. arXiv:2111.00364 [cs.LG]

[53] Ya Xu et al. 2015. From infrastructure to culture: A/B testing challenges in large scale social networks. In *KDD*. 2227–2236.

[54] Dragomir Yankov, Pavel Berkhin, and Lihong Li. 2015. Evaluation of explore-exploit policies in multi-result ranking systems. *arXiv:1504.07662* (2015).

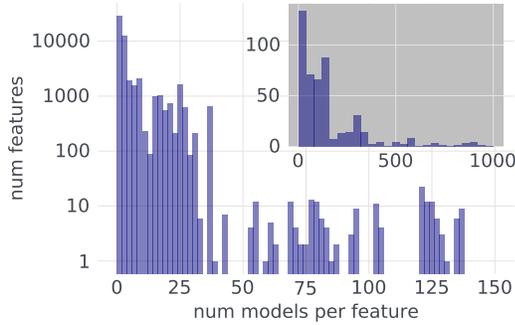[55] Zhe Zhao et al. 2019. Recommending what video to watch next: a multitask ranking system. In *RecSys '19*. 43–51.

Fig. 6. Histograms showing how many Looper models use a given feature (outer) and features per model (inner).



Fig. 8. The three phases of resource optimization via *feature reaping*. Resource consumption is decreased without adverse product impacts.

## A RESOURCE UTILIZATION AND ITS OPTIMIZATION

Smart strategies tend to provide significant benefits but sometimes need serious computational resources, Therefore product deployment requires judicious resource management. Looper is deployed in numerous and diverse applications at Meta, some of which optimize performance of other systems and some enhance functionality (Figure 7). Such economies-of-scale infrastructure enables resource reuse and load-balancing. Figure 5 shows that different use cases exhibit different model-lifecycle bottlenecks, with *feature extraction* drawing the largest share of resources for demanding use cases.[1] Based on this trend, we developed a *feature reaping* optimization that removes unimportant features. This optimization estimates the importance of individual features, removes features in groups, then checks the results by training a reduced model and evaluating it. When we deployed feature reaping in production, it and provided overall 11% resource-cost savings (10-30% per use case) with no adverse product impacts. Figure 8 illustrates the resource savings (28%) for one use case via feature reaping. It distinguishes 3 stages that experience constant decision traffic: (*i*) offline model training that runs feature reaping and trains a reduced model, (*ii*) the online stage that uses the old and the revised (20% traffic only) model to evaluate product impact, and (*iii*) the production stage that uses only the new model whose performance had been validated by online evaluation (which took 2 weeks to collect statistically significant results). Sharing such AutoML optimizations across multiple use cases makes our platform competitive with specialized platforms.
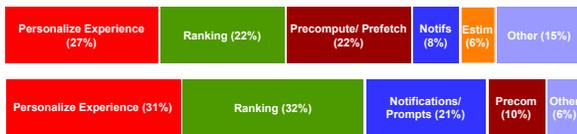


Fig. 7. Product adoption of smart strategies by use case count (top) and by resource consumption (bottom). Resource consumption during training is correlated with data amount, but resource consumption at inference reflects decision rates in applications, the number of features used by models and the presence of synthetic/engineered features.

---

[1]We note that Looper deploys moderate-complexity models with *diverse metadata features*, whereas advanced deep learning models with *homogeneous* image pixels, word embeddings, etc may exhibit different trends.

## B A SURVEY OF PLATFORM ADOPTERS

The article "Why Machine Learning Strategies Fail" [21] lists common barriers to entry:

- lacking a business case,
- lacking data,
- lacking ML talent,
- lacking sufficient in-house ML expertise for outsourcing,
- failing to evaluate an ML strategy.

When talking to prospective clients, we advise against using Looper when there is no need for an end-to-end platform. For example, in a Kaggle-like environment with well-defined data, the goal is to train a model to optimize closed-form objectives for model performance. A second example is tasks without a clear product metric, such as building latent-space embeddings and other self-supervised tasks. Yet another reason not to use Looper is the efficiency of special-case platforms, e.g., for ranking and image-processing.

To clarify when Looper is relevant and to clarify the adoption process of smart strategies, we interviewed product teams at Meta that adopted our platform and saw product impacts (Figure 7). All the teams had tried heuristic approaches but with poor results, hence their focus on ML. Simple heuristics proved insufficient for user bases spanning multiple countries with distinct demographic and usage patterns. The following challenges were highlighted: (1) manually optimizing parameters in large search spaces, (2) figuring out the correct rules to make heuristics effective, (3) trading off multiple objectives, (4) updating heuristic logic quickly, especially in on-device code.

The spectrum of ML expertise varied across product teams from beginners to experienced ML engineers, and only 15% of teams using our platform include ML engineers. For teams without production ML experience, an easy-to-use ML platform is often the deciding factor for ML adoption, and ML investment continues upon evidence of utility. An engineer mentioned that a lower-level ML system had a confusing development flow and unwieldy debugging. They were unable to set up recurring model training and publishing. Looper abstracts away concerns about SW upgrades, logging, monitoring, etc behind high-level services and unlocks hefty productivity savings.

For experienced ML engineers, a smart-strategies platform improves productivity by automating repetitive time-consuming work:

writing database queries, implementing data pipelines, setting up monitoring and alerts. Compared to narrow-focus systems, it helps product developers launch more ML use cases. An engineer shared prior experience writing custom queries for features and labels, and manually setting up pipelines for recurring training and model publishing without an easy way to monitor model performance and issue emergency alerts. Some prospective clients who evaluated our platform chose other ML platforms within our company or stayed with their custom-designed infrastructure. They missed batched offline prediction with mega-sized data and needed exceptional performance possible only with custom ML models. These issues can be addressed with additional platform development efforts.

Successful platform adopters configured initial ML models in two days and started collecting training data. Training the model using product feedback and revising it over 1-2 weeks enabled online product experiments that take 2-4 weeks. HTE analysis and impact optimization take 1-3 weeks. Product launch can occur 1-3 months after initial data collection. Among platform adopters, experienced engineers aware of ML-related technical debt and risks [2, 21, 40, 45, 46] appreciated the built-in support for recurring training, model publishing, data visualization, as well as monitoring label and feature distributions over time with data-drift alerts. Also noted was the canarying mechanism for new models (Section 3.2). Surprisingly important was helping adopters share model insights (such as feature importance) with leadership and product managers. Among possible improvements, adopters mentioned development velocity.

## C   IMPROVING PRODUCT-DEVELOPMENT VELOCITY

As a platform used primarily by product engineers, Looper aims to reduce friction while users modify parameters and model configurations. This is accomplished in two ways:

(1) a GUI to easily edit all Looper parameters, and
(2) automatic optimization of configuration parameters.

Every aspect of configuration is represented in the UI. Each use case has a landing page to manage all related models, data versions, and experiments. All parameters including features, labels, and current production models are displayed in an editable form for users to modify and immediately push changes to production. Experiments can be launched from the use case landing page and the UI allows users to select baseline and experiment models to immediately compare them using integrated A/B tests.

All parameters are stored in configuration and this enables automatic optimization of parameters. Using the integrated experiment optimization system [9], several values for a parameter are optimally selected to immediately test in our integrated A/B experimentation platform. These values can then be tested against product metrics. Since most parameters required for experimentation are already stored in strategy blueprints, including eventual support for product metrics, it is possible to automatically improve various parameter values using the experiment optimization system and A/B test to check if the config optimization benefits product, system, and resource usage metrics without affecting production models and without use case owner input.

## D   DEVELOPMENT EFFORT

The Looper platform described in this paper was developed at Meta over several years. It uses several types of software infrastructure, such as databases, horizontal ML platforms, reusable ML models and frameworks, product metrics, and support for product A/B testing. The overall design was revised to better adapt to the needs of applications. On the other hand, a team of ten experienced software engineers should be able to implement our core platform design in half a year using relevant open-source and/or company infrastructure. In such development, it is important to focus on representative product use cases and guide software development within a well-defined scope. Avoiding common problems, rather than developing comprehensive solutions, can reduce time to first application. In particular, the complete chain of custody of data in Looper helps avoid or reduce many common problems with data quality, such as delayed and missing data, mismatches between training and testing, etc.