
A Bandit Approach to Posterior Dialog Orchestration Under a Budget

Sohini Upadhyay, Mayank Agarwal, Djallel Bounneffouf, and Yasaman Khazaeni
IBM Research AI
{firstname.lastname}@ibm.com

Abstract

Building multi-domain AI agents is a challenging task and an open problem in the area of AI. Within the domain of dialog, the ability to orchestrate multiple independently trained dialog agents, or skills, to create a unified system is of particular significance. In this work, we study the task of online posterior dialog orchestration, where we define posterior orchestration as the task of selecting a subset of skills which most appropriately answer a user input using features extracted from both the user input and the individual skills. To account for the various costs associated with extracting skill features, we consider online posterior orchestration under a skill execution budget. We formalize this setting as Context Attentive Bandit with Observations (CABO), a variant of context attentive bandits, and evaluate it on simulated non-conversational and proprietary conversational datasets.

1 Introduction

Serverless execution enables scalable and modular deployment of models for contemporary applications, including AI agents. In the context of dialog systems, such modular design entails connecting multiple dialogue agents or "skills" - each trained independently on different or overlapping tasks - to form a unified system with capabilities of each of its constituent skills. The orchestrator or control module for such a system can either be a deterministic module employing "If this then that" (IFTTT) logic or more complex functional programming frameworks such as Amazon Lambda etc, or could in itself be a learnable model employing either supervised or reinforcement learning approaches.

This is a fairly common scenario in contemporary personal home assistant devices such as Amazon Alexa and Google Home, where developers have the ability to integrate their own independently developed skills with the assistant's core infrastructure. Here, the assistant itself is responsible for invoking skills in response to user input. Invocation of these skills falls into two categories: explicit invocation and implicit invocation. Explicit invocation occurs when the user explicitly specifies the name of the skill they are interested in interacting with. This requires the user to specify the name of the skill along with a set of pre-defined invocation phrases that trigger the skill [1], and is an application of the IFTTT logic. Implicit invocation on the other hand, does not provide the assistant with the name of the skill the user is interested in interacting with, and requires the assistant to understand the users query along with the available skills capabilities to select the most appropriate skill to respond with [1, 2]. Implicit invocation has a clear advantage in facilitating more natural conversation and removes the knowledge barrier of skill naming and understanding that explicit invocation requires.

Dialog orchestration models that use implicit invocation tend to follow the a-priori approach or the posterior approach. A-priori orchestrator models are built exclusively using features known prior to executing any skills whereas posterior models execute skills to extract supplemental features. By this definition, it is clear that posterior approaches ought to match or beat comparable a-priori methods

as they use a superset of the features used in the a-priori approach. Most recent work has focused on the posterior approach, including submissions to the Alexa prize competition [3, 4]. Both of these approaches have leveraged supervised learning techniques, which necessitate training data and regular updating if deployed live.

Online orchestration models could remove this hurdle, enabling a cold start deployment. Online orchestration also does not require a fixed label space, allowing new skills to be added to the agent in a seamless way. While a multitude of reinforcement learning models are viable orchestration candidates, we investigate the use of contextual bandits for the task. The contextual bandit problem is a variant of the extensively studied multi-armed bandit problem [5, 6, 7], where at each iteration, before choosing an arm, the agent observes an N -dimensional *context*, or *feature vector*, and uses it to predict the next best arm to play [8, 9, 10, 11]. Every time an arm is played, a reward value is observed. Over time, the agent’s aim is to collect enough information about the relationship between the context vectors and rewards, so that it can predict the next best arm to play by looking at the corresponding context [8, 11].

Posterior orchestration, online or otherwise, is not without its challenges. Recall that in posterior dialog orchestration, a user’s query is often directed to a number of domain specific skills and the best response is returned. In this case, the pre-execution features, i.e features extracted from the query, can be immediately observed, but the set of features or responses from skills, the post-execution features, cannot. For multi-purpose dialog systems, like personal home assistants, executing and retrieving features or responses from every skill can be computationally expensive or intractable, with the potential to cause a poor user experience. Moreover, executing skills in some use cases may necessitate api requests associated with actual costs. Thus while posterior dialog orchestration models are in many ways conceptually preferable to a-priori approaches, in practice they are associated with an often unaccounted cost. The challenge here is introducing a budget on the number of post-execution features that can be extracted. Some existing supervised posterior orchestration methods recognize this challenge and avoid retrieving all post-execution features. As an example, Kim et. al. [12] present a set of efficient and scalable neural shortlisting-reranking models for personal assistants. The shortlisting stage efficiently trims all the skills down to a list of top-k candidates, and the reranking stage performs a list-wise reranking of the initial top-k skills with additional contextual information. Beyond the lack of cold-start support inherent to all supervised approaches, the amount of data necessary to effectively train this kind of neural model based method is a limiting factor for low-data use cases. Given that online orchestration avoids these pitfalls, we develop a novel bandit algorithm that handles this challenge of limited access to post-execution features.

The goal of our research is to build a dialog orchestration framework which can utilize query and user features along with the conversational context to route the dialog in a multi-skill system. Overall, the main contributions of this paper include (1) presenting an online approach to dialog orchestration, (2) a new variant of the context attentive bandit problem, motivated by limitations of posterior dialog orchestration, and (3) an empirical evaluation demonstrating the advantages of our proposed method over a range of datasets and settings.

2 Background

The contextual bandit problem has been extensively studied in the past, and a variety of solutions have been proposed. In LINUCB [13, 14, 15] and in Contextual Thompson Sampling (CTS) [11], a linear dependency is assumed between the expected reward given the context and an action is taken after observing this context; the representation space is modeled using a set of linear predictors. However, the context is assumed to be *fully observable*, which is not the case in this work.

Motivated by dimensionality reduction tasks, Abbasi-Yadkori et. al. [16] studied a sparse variant of stochastic linear bandits, where only a relatively small and unknown subset of features is relevant to a multivariate function optimization. Similarly, Carpentier & Munos [17] also considered high-dimensional stochastic linear bandits with sparsity, where s components are assumed to be non-zero, and where the dimension N of the context is larger than the sampling budget n . In Bastani & Bayati [18] consider a multi-arm bandit (MAB) problem with high-dimensional covariates, and a new efficient bandit algorithm based on the LASSO estimator is presented. Regret analysis is performed, demonstrating that the proposed algorithm achieves near-optimal performance in comparison to

an oracle that knows all the problem parameters. Still, all above work, unlike ours, assumes full observability of the context variables, which is not the case in many important applications.

Finally, Bouneffouf et. al. [19] developed the idea of context attentive bandits - a case of the contextual bandit problem, referred to as contextual bandit with restricted context (CBRC), where observing the whole feature vector at each iteration is impossible, and the agent can only request to see some limited number of those features; the upper bound (budget) on the feature subset is fixed for all iterations, but within this budget, the agent can choose any feature subset of said size. However, in the posterior dialog orchestration application, while the full context may be too costly or impossible to see, some partial observation of the context, e.g. query or user features, can be known to an agent initially, along with the ability to observe unknown context features, up to a certain limit, as in CBRC.

Motivated by the limitations of posterior dialog orchestration, we extend the context attentive bandit to a special case which we call the *Context Attentive Bandit with Observations (CABO)*. In the CABO setting, observing the full context vector at each iteration is impossible, but a small subset of context features, is observable and a fixed number of the unobserved features within a budget can be revealed. The goal here is to leverage the observable features to select the best unknown feature subset at each iteration to maximize overall reward.

3 Problem Setting

We begin by formally defining concepts our novel bandit problem setting builds upon, such as contextual bandit and contextual combinatorial bandit.

The Contextual Bandit Problem. Following Langford & Zhang [8], this problem is defined as follows. At each time point (iteration) $t \in \{1, \dots, T\}$, an agent is presented with a *context* (feature vector) $\mathbf{c}(t) \in \mathbf{R}^N$ before choosing an arm $k \in A = \{1, \dots, K\}$. We denote by $C = \{C_1, \dots, C_N\}$ the set of features (variables) defining the context. Let $\mathbf{r}(t) = (r_1(t), \dots, r_K(t))$ denote a reward vector, where $r_k(t) \in [0, 1]$ is a reward at time t associated with the arm $k \in A$. Herein, we will primarily focus on the Bernoulli bandit with binary reward, i.e. $r_k(t) \in \{0, 1\}$. Let $\pi : \mathbf{R}^N \rightarrow A$ denote a policy, mapping a context $\mathbf{c}(t) \in \mathbf{R}^N$ into an action $k \in A$. We assume some probability distribution $P_c(c)$ over the contexts in C , and a distribution of the reward, given the context and the action taken in that context. We assume that the expected reward (with respect to the distribution $P_r(r|\mathbf{c}, k)$) is a linear function of the context, i.e. $E[r_k(t)|\mathbf{c}(t)] = \mu_k^T \mathbf{c}(t)$, where μ_k is an unknown weight vector associated with the arm k ; the agent’s objective is to learn μ_k from the data so it can optimize its cumulative reward over time.

Contextual Combinatorial Bandit. Our feature subset selection approach builds upon the *Contextual Combinatorial Bandit (CCB)* problem [20], specified as follows. Each arm $k \in \{1, \dots, K\}$ is associated with the corresponding variable $x_k(t) \in R$ indicating the reward obtained when choosing the k -th arm at time t , for $t > 1$. In the contextual combinatorial bandit setting, the agent sequentially observes a context \mathbf{c} , selects a subset of arms $M \in S$, from a constrained set of arm subsets $S \subseteq P(K)$, where $P(K)$ is the power-set of K , and observes a reward $r_M(t) = h(x_k(t)), k \in M$ associated with the selected subset of arms. Here we define the reward function $h(\cdot)$ used to compute $r_M(t)$ as a sum of the outcomes of the arms in M , i.e. $r_M(t) = \sum_{k \in M} x_k(t)$, although one can also use nonlinear rewards. The objective of the CCB algorithm is to maximize the reward over time. We consider here a stochastic model, where the expectation of $x_i(t)$ observed for an arm k is a linear function of the context, i.e. $E[x_i(t)|\mathbf{c}(t)] = \mu_i^T \mathbf{c}(t)$, where μ_i is an unknown weight vector (to be learned from the data) associated with the arm i . The distributions can be different for each arm. The global rewards $r_M(t)$ are also random variables, independent and distributed according to some unknown distribution with some expectation μ^M .

3.1 CABO: Context Attentive Bandit with Observations

Building off the contextual bandit and contextual combinatorial bandit problems, we formally define a novel type of bandit problem, called *Context Attentive Bandit with Observations (CABO)*.

As mentioned above, $\mathbf{c}(t) \in \mathbf{R}^N$ will denote a vector of values assigned to an (ordered) set of random context variables, or features, $C = \{C_1, \dots, C_N\}$, at time t . Let $C^D \subseteq C$, $|C^D| = D$, $0 < D \leq N$, denote a subset of features of size D , and let $\mathbf{c}^D(t) \in S_{C^D}$ denote a vector from a D -subspace of

\mathbf{R}^N , denoted $S_{C^D} \subseteq \mathbf{R}^N$, which is defined as a subspace containing all sparse vectors with features (coordinates) outside of the subset C^D set to zero.

We assume that at each time point t the environment generates a feature vector $\mathbf{c}(t) \in \mathbf{R}^N$ which the agent cannot observe fully. However, unlike the previously introduced CBRC setting [19], the agent has now a *partial observation* of the context, i.e. it can see a small subset of *observed* features $C^O \subset C$, where $O \ll N$. Given these observations \mathbf{c}^O , the agent is allowed to request more features to observe (similar to CBRC setting), up to D (desired) features in total, including the initial set $C^O \subset C^D$, where C^D denotes final set of D observed features. Assuming that the unobserved features are all of the same fixed cost, there is a budget of $U = D - O$ features imposed on the agent. The goal of the agent is to maximize its total reward over time via (1) the optimal choice of the additional observations, given the initial ones, and (2) the optimal choice of a subsequent action $k \in A$ based on the resulting extended observation.

Let us now formally define the set of all policies, i.e. possible mappings from agent’s observations to its actions restricted to the proposed problem setting, as the set of the compound functions

$$\begin{aligned} \Pi_o^D = \cup_{C^O \subseteq C} \{ \pi : S_{C^O} \rightarrow A, \text{ s.t. } \pi(\mathbf{c}^O) = \hat{\pi}(\mathbf{c}^D), \\ \mathbf{c}^D = g(C^D), C^D = h(C^O, \mathbf{c}^O) \}, \end{aligned} \quad (1)$$

where

- $g : P_D(N) \rightarrow S_D(\mathbf{R}^N)$ is a function mapping a given subset of features $C^D \in P_D(N)$, $P_D(N)$ denoting the set of all subsets of $\{1, \dots, N\}$ of size D , to a vector $\mathbf{c}^D \in S_D(\mathbf{R}^N)$, $S_D(\mathbf{R}^N)$ denoting the set of all d -subspaces S_{C^D} of \mathbf{R}^N , each defined for a corresponding subset C^D of features;
- $h : P_O(N) \rightarrow P_D(N)$ maps the initial set of observed features C^O to the extended set of features to be observed, C^D , $C^O \subset C^D$;
- $\hat{\pi} : S_{C^D} \rightarrow A$ is a function mapping the observed extended feature subset \mathbf{c}^D into an action (a.k.a. bandit’s arm) $k(t) \in A$, which results into a reward $r_k(t)$.

The objective of a contextual bandit algorithm is to find an optimal policy $\pi \in \Pi_o^D$, over T iterations or time points, so that the total reward is maximized.

4 Methodology

4.1 CATSO: Context Attentive Thompson Sampling with Observations

We propose a novel method for solving the CABO problem, which we name *Context Attentive Thompson Sampling with Observations (CATSO)*, and summarize it in Algorithm 1. The combinatorial task of selecting the best subset of features is treated as a contextual combinatorial bandit (CCB) problem [20], and the subsequent decision-making (action selection) task as a contextual bandit problem solved by Contextual Thompson Sampling (CTS) [11], respectively.

The algorithm takes the total number of features N , initially observed number of features O , and the total desired number of features to observe D , as inputs. We use $U = D - O$ to denote our budget, the number of unobserved features to reveal. We will use several stages, up to S , to reveal U features. When $S = 1$, the observed features O are used to select all U features as a set, whereas when $S = U$, the set of O features is updated incrementally and used to select each of the U additional features one at a time. The algorithm also requires hyperparameter v , the exploration parameter used in Thompson Sampling.

The algorithm iterates over T steps, where at each iteration t , the values $\mathbf{c}^O(t)$ of features in the original observed subset C^O are observed first. The current set of already observed features, X^t , and the corresponding observed context, $\mathbf{x}(t)$, is maintained over all stages, and are initialized to C^O and \mathbf{c}^O respectively. At each iteration t , the vector parameter θ_i is sampled from the corresponding multivariate Gaussian distribution (step 10) for each feature i not yet observed so far, to estimate $\hat{\theta}_i$. Thereafter, at each stage, the best subset of features are selected, $C^u \subseteq C/X^t$, such that $C^u = \arg \max_{C' \subseteq C/X^t, |C'|=u} \sum_{i \in C'} \mathbf{x}(t)^\top \theta_i$ where $u = U/S$ is the number of unknown features to explore at each stage.

Algorithm 1 Context Attentive Thompson Sampling with Observations

- 1: **Require:** Total number of features N ; initially observed number of features O ; the set of those features and their values, C^O and c^O , respectively; the total desired number of features to observe D , over S stages; the exploration parameter v , and the function $\lambda(t)$ computed differently for stationary and nonstationary cases.
 - 2: **Initialize:** $\forall k \in \{1, \dots, K\}, A_k = I_K, g_k = 0_K, \hat{\mu}_k = 0_K$, and $\forall i \in \{1, \dots, N\}, B_i = I_N, z_i = 0_N, \hat{\theta}_i = 0_N$.
 - 3: $U = D - O, u = U/S$
 - 4: **Foreach** $t \in 1, 2, \dots, T$ **do**
 - 5: observe $c^O(t)$, given feature subset C^O
 - 6: $X^t = C^O, \mathbf{x}(t) = c^O$
 - 7: **Foreach** stage $s = 1, 2, \dots, S$ **do**
 - 8: **Foreach** context feature $i = 1, \dots, N$ **do**
 - 9: **if** $i \notin X^t$ **then**
 - 10: Sample θ_i from $\mathcal{N}(\hat{\theta}_i, v^2 B_i^{-1})$
 - 11: **End if**
 - 12: **End do**
 - 13: Select $C^u(t) = \underset{C' \subseteq C/X^t, |C'|=u}{\operatorname{argmax}} \sum_{i \in C'} \mathbf{x}(t)^\top \theta_i$
 - 14: $X^t = X^t + C^u(t)$
 - 15: observe values $\mathbf{x}(t)$ of feature subset X^t
 - 16: **End do**
 - 17: **Foreach** arm $k = 1, \dots, K$ **do**
 - 18: Sample μ_k from $\mathcal{N}(\hat{\mu}_k, v^2 A_k^{-1})$ distribution.
 - 19: **End do**
 - 20: Select arm $k(t) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \mathbf{x}(t)^\top \mu_k$
 - 21: Observe $r_k(t)$
 - 22: $A_k = A_k + \mathbf{x}(t)\mathbf{x}(t)^\top$
 - 23: $g_k = g_k + \mathbf{x}(t)r_k(t)$
 - 24: $\hat{\mu}_k = A_k^{-1}g_k$
 - 25: **Foreach** $i \in X^t \setminus C^O$
 - 26: $B_i = \lambda(t)B_i + \mathbf{x}(t)\mathbf{x}(t)^\top$
 - 27: $z_i = z_i + \mathbf{x}(t)r_k(t)$
 - 28: $\hat{\theta}_i = \lambda(t)B_i^{-1}z_i$
 - 29: **End do**
 - 30: **End do**
-

Once a subset of features is selected using the contextual combinatorial bandit approach, the algorithm switches to the contextual bandit setting to choose an arm based on the context consisting now of a subset of features (steps 17-24).

We assume that the expected reward is a linear function of a restricted context,

$$E[r_k(t)|\mathbf{x}(t)] = \mu_k^\top \mathbf{x}(t)$$

We assume that reward $r_k(t)$ for choosing arm k at time t follows a parametric likelihood function $P(r(t)|\mu_k)$, and that the posterior distribution at time $t + 1$, $P(\mu|r(t)) \propto P(r(t)|\mu)P(\mu)$, is given by a multivariate Gaussian distribution $\mathcal{N}(\hat{\mu}_k(t + 1), v^2 A_k(t + 1)^{-1})$ where $A_k(t) = I_N + \sum_{\tau=1}^{t-1} c(\tau)c(\tau)^\top$ with N the size of the context vectors c , and $\hat{\mu}_k = A_k(t)^{-1}(\sum_{\tau=1}^{t-1} c(\tau)c(\tau))$.

At each time point t , and for each arm, a k -dimensional μ_k is sampled from $\mathcal{N}(\hat{\mu}_k(t), v^2 A_k(t)^{-1})$, an arm is chosen such that $\mathbf{x}(t)^\top \mu_k$ is maximized (step 20 in the algorithm), a reward $r_k(t)$ is obtained for choosing an arm k , and finally the relevant parameters are updated.

4.1.1 CATSO in Nonstationary Setting

Practical posterior dialog orchestration applications motivate the need to consider the possibility of nonstationary unobserved context features. In posterior dialog orchestration, we assume each domain

specific skill outputs features pertaining to their query response. In some use cases, each skill could be independently updated at any time, changing these features. As a result, similar queries, which would likely define the observable context C^O , can elicit vastly different distributions of response features, the unknown context, over time. The main problem with any stationary algorithm is that it gives equal weight to its history. In a nonstationary environment, if there is no specific assumption about how the environment will change, a simple idea is to use a weighting function to lessen the effect of the past on current decisions. Since we are using CTS as our base model and it uses ridge regression, implementation of weighting instances is straightforward. We propose assigning decaying weights to the past examples in the ridge regression. The same kind of weights are also applied in the calculation of the confidence width. Following the notation from Algorithm 1, in this case $\lambda(t)$ represents the decay parameter. In order to compute the optimal $\lambda(t)$ value, we use GP-UCB algorithm [21], which is an algorithm that solves the multi-armed bandit problem in continuous space. Computing $\lambda(t)$ is done via the following decision rule:

$$\lambda(t) = \operatorname{argmax}_{\lambda \in D} [\mu_{t-1}(\lambda) + \alpha^{1/2} \sigma_{t-1}(\lambda)]$$

with α as the GP-UCB exploration parameter, μ_{t-1} mean reward, and σ_{t-1} the uncertainty. The GP-UCB algorithm is initialized with the search space $\lambda \in [0, 1]$ and at each iteration uses the above equation to calculate a different $\lambda(t)$, which is then used by CATSO. More detail on the GP-UCB algorithm can be found in [21].

5 Experiments

We assess Context Attentive Thompson Sampling with Observations (CATSO) with respect to the current state of the art for context attentive bandits, Thompson Sampling with Restricted Context (TSRC). TSRC solves the contextual bandit with restricted context problem (CBRC) discussed prior, which selects a set of unknown features at each event while assuming no observable features exist initially. For a total number of features N , we refer to the O observed features as the known context and the $N - O$ unobserved context features as the unknown context. In our use of the TSRC algorithm, at each iteration, the known context is observed, the TSRC decision mechanism independently chooses U unknown context features to reveal, and Contextual Thompson Sampling (CTS) is invoked. Empirical evaluation of CATSO and TSRC was performed on publicly available classification datasets and on a propriety corporate dialog orchestration dataset.

Publicly available Coverttype¹ and CNAE-9¹ were featured in the original TSRC paper and Warfarin [22] is a historically popular dataset for evaluating bandit methods. The details of these datasets are summarized in Table 1.

For the stationary setting, we randomly fix 10% of the context feature space of each dataset to be known at the onset and explore a subset of U unknown features. For CATSO, we fix $\lambda(t) = 1$ to reflect the stationary setting and choose $S = 1$. For the nonstationary setting, we simulate nonstationarity in the unknown feature space by duplicating each dataset, randomly fixing the known context in the same manner as above, and shuffling the unknown feature set - label pairs. Then we stochastically replace events in the original dataset with their shuffled counterparts, with the probability of replacement increasing uniformly with each additional event. For this nonstationary setting, which we refer to as NCATSO, we again fix $S = 1$, but use $\lambda(t)$ as defined by the GP-UCB algorithm. We compare NCATSO to Weighted TSRC (WTSRC), the nonstationary version of TSRC also developed by Bouneffouf et al. [19]. WTSRC makes updates to its feature selection model based only on recent events, where recent events are defined by a time period, or "window" w . We choose $w = 100$ for WTSRC. We report the total average reward across a range of U corresponding to various percentages of N for each algorithm in each setting in Table 2. The results in Table 2 are promising, with our methodologies outperforming the state of the art in the majority of cases across both settings. The most notable exception is where CATSO sometimes outperforms and other times nearly matches TSRC performance on the CNAE-9 dataset. This outcome is somewhat expected, for in the original work on TSRC [19], the mean error rate of TSRC was only 0.03% lower than randomly fixing a subset of unknown features to reveal for each event on CNAE-9. This suggests that the operating

Table 1: Datasets

Datasets	Instances	Features	Classes
Coverttype	500 000	95	7
CNAE-9	1080	856	9
Warfarin	5528	93	3

¹<https://archive.ics.uci.edu/ml/datasets.html>

Table 2: Total average reward, $O = 10\%$

(a) Stationary setting				(b) Nonstationary setting			
	Warfarin				Warfarin		
U	20%	40%	60%	U	20%	40%	60%
TSRC	53.28 ± 1.08	57.60 ± 1.16	59.87 ± 0.69	WTSRC	55.83 ± 0.55	58.00 ± 0.83	59.85 ± 0.60
CATSO	53.65 ± 1.21	58.55 ± 0.67	60.40 ± 0.74	NCATSO	59.47 ± 2.89	59.34 ± 2.04	63.26 ± 0.75
	Covertime				Covertime		
U	20%	40%	60%	U	20%	40%	60%
TSRC	54.64 ± 1.87	63.35 ± 1.87	69.59 ± 1.72	WTSRC	50.26 ± 1.58	58.99 ± 1.81	64.91 ± 1.38
CATSO	65.57 ± 2.17	72.58 ± 2.36	78.58 ± 2.35	NCATSO	48.50 ± 1.05	68.17 ± 3.14	83.78 ± 5.51
	CNAE-9				CNAE-9		
U	20%	40%	60%	U	20%	40%	60%
TSRC	33.57 ± 2.43	38.62 ± 1.68	42.05 ± 2.14	WTSRC	19.91 ± 2.67	30.86 ± 2.92	36.01 ± 2.88
CATSO	29.84 ± 1.82	39.10 ± 1.41	40.52 ± 1.42	NCATSO	30.88 ± 0.96	34.91 ± 1.93	42.04 ± 1.52

premise of TSRC, that some features are more predictive of reward than others, does not hold on this dataset. On top of this assumption, CATSO also assumes that there exist relationships between the known and unknown context features, likely causing a small compounding of error.

Proprietary corporate dialog application *Customer Assistant* orchestrates 9 domain specific skills which we arbitrarily denote as $Skill_1, \dots, Skill_9$ in the discussion that follows. In this application, example skills lie in the domains of payroll, compensation, travel, health benefits, and so on. Each skill is designed with a multi-turn conversation dialog tree. In addition to a textual response to a user query, the skills orchestrated by *Customer Assistant* also return the following features: an *intent*, a short string descriptor that categorizes the perceived intent of the query, and a *confidence*, a real value between 0 and 1 indicating how confident a skill is that its response is relevant to the query. Skills have multiple intents associated with them. The orchestrator uses all the features associated with the query and the candidate responses from all the skills to choose which skill should carry the conversation at a given event.

We accessed the training data for each skill to find example queries *Customer Assistant* has valid responses for, amounting to 28,412 queries total. Accordingly, we denote the correct class for a query to be the skill it was an example for. For 127 queries common to more than one skill’s training data, one of the skills was randomly assigned as the correct class. We pose each query to all of the skills to extract the associated intents and confidences, and add noise sampled from a $\mathcal{N}(-0.1, 0.05)$ distribution to all of the confidences to avoid built-in biases. We encode each query by averaging 50 dimensional GloVe word embeddings [23] for each word in each query and for each skill we create a feature set consisting of its confidence and a one-hot encoding of its intent. The skill feature set size for $Skill_1, \dots, Skill_9$ are 181, 9, 4, 7, 6, 27, 110, 297, and 30 respectively. We concatenate the query features and all of the skill features to form a 721 dimensional context feature vector for each event in this dataset. In contrast to the publicly available datasets, here there is no need for simulation of the known and unknown contexts; in a live setting the query features are immediately calculable or known, whereas the confidence and intent necessary to build a skill’s feature set are unknown until a skill is executed. Because the confidence and intent for a skill are both accessible post execution, we reveal them together. We accommodate this by slightly modifying the objective of CATSO to reveal U unknown skill feature sets instead of U unknown individual features for each event.

We perform a deeper analysis of the *Customer Assistant* dataset, examining the case where $S = U$. Recall that when $S = 1$, the known context, in this case the query features, is used to select all U additional context features sets at once, whereas when $S = U$, the known context grows and is used to select each of the U additional context feature sets incrementally. Maintaining $\lambda(t) = 1$, for the stationary case we denote these two cases of the CATSO algorithm as CATSO-1 and CATSO-U respectively and report their performance across various U , the number of unknown skill feature sets revealed. Note that when all 9 skill feature sets are revealed, the CATSO and TSRC methods all reduce to simple Contextual Thompson Sampling (CTS) with the full feature set. Similarly, when 0 skill feature sets are revealed, the methods all reduce to CTS with a sparsely represented context of the query features. CTS suffers from this sparsity so we also consider a case we call CTS-query, CTS where the context is exclusively the query features. CTS-query is thus an a-priori online approach to dialog orchestration that completely ignores the existence of post-execution features. The results for the stationary case are summarized in Figure 1. CATSO-U appears to slightly outperform CATSO-1

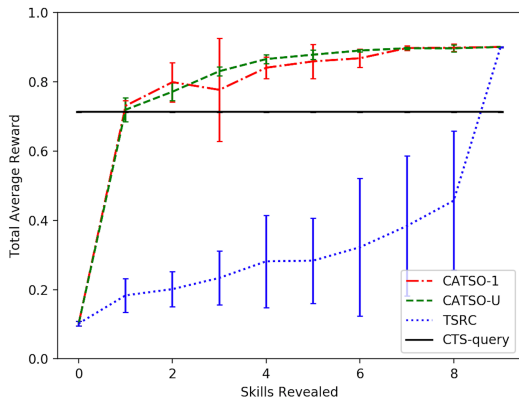


Figure 1: Stationary Setting - Customer Assistant with 9 Skills

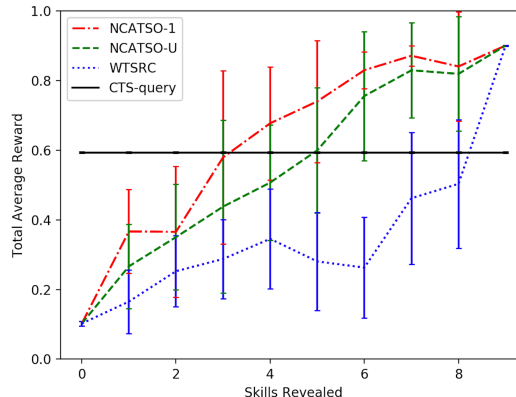


Figure 2: Nonstationary Setting - Customer Assistant with 9 Skills

across all U tested and both methods outperform TSRC by a large margin. Also notice that our posterior methods CATSO-1 and CATSO-U outperform the a-priori method CTS-query even under very small post-execution feature budgets, as low as 2 skill feature sets.

For the nonstationary case we simulate nonstationarity in the same manner as the publicly available datasets, except using the natural partition of the query features as the known context and the skill feature sets as the unknown context instead of simulated percentages. We use the GP-UCB algorithm for $\lambda(t)$, refer to the $S = 1$ and $S = U$ cases as NCATSO-1 and NCATSO-U, and illustrate their performance alongside WTSRC and CTS-query in Figure 2. Here we observe that NCATSO-1 slightly outperforms NCATSO-U, and both outperform the WTSRC baseline. Notice that posterior approach NCATSO-1 outperforms CTS-query, the a-priori approach, when approximately 3 or more skill feature sets are revealed.

6 Conclusions and Future Work

In this paper we consider how to address the challenges of posterior dialog orchestration using an online approach. We formulated CABO, a new variant of context attentive bandits motivated by practical budgets on skill execution and demonstrate that our new bandit algorithm beats the existing state of the art context attentive bandit algorithm on simulated (non-dialog) and dialog datasets across stationary and nonstationary settings.

Customer Assistant has now been deployed to over 100,000 users with a thumbs up/down feature that allows users to provide individual feedback to each of the responses. The system also allows Subject Matter Experts (SMEs) to provide explicit labels to each event, enabling human evaluation of all the models.

Theoretical regret bounds on the proposed algorithm will follow in a more theory focused work. Our current algorithm treats multi-skill dialog orchestration as a single class problem, where after each query, only one skill response is returned to the user. However, in many use cases, multiple responses ought to be returned to the user. We would like to shift our algorithm to the multiclass setting, perhaps by using the contextual combinatorial bandit approach in the arm selection process in addition to its current role in the unknown context feature selection process. Also, our current formulation assumes that all of the unobserved features are of the same cost, and thus the budget on cost is equivalent to a budget on the number of features. We plan on expanding this notion of budget to accommodate settings where unobserved features have different costs. Other directions for future work include using non-bandit algorithms in the context feature selection stage and exploring nonstationarity in the known context space.

References

- [1] Google. Overview | actions on google | google developers. <https://developers.google.com/actions/discovery/>, 2018.

- [2] Amazon. Understand how users invoke custom skills. <https://developer.amazon.com/docs/custom-skills/understanding-how-users-invoke-custom-skills.html>, 2018.
- [3] Ioannis Papaioannou, Amanda Cercas Curry, Jose L Part, Igor Shalyminov, Xinnuo Xu, and Yanchao Yu. Alana : Social Dialogue using an Ensemble Model and a Ranker trained on User Feedback. *1st Proceedings of Alexa Prize*, pages 1–10, 2017.
- [4] Oluwatosin Adewale, Alex Beatson, Davit Buniatyan, Jason Ge, Mikhail Khodak, Holden Lee, Niranjani Prasad, Nikunj Saunshi, Ari Seff, Karan Singh, Daniel Suo, Cyril Zhang, and Sanjeev Arora. Pixie : A Social Chatbot. *Alexa Price Proceedings 2017*, pages 1–10, 2017.
- [5] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- [6] John C Gittins. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(2):148–164, 1979.
- [7] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [8] John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. In *Advances in neural information processing systems*, pages 817–824, 2008.
- [9] Deepak Agarwal, Bee-Chung Chen, and Pradheep Elango. Explore/exploit schemes for web content optimization. In *2009 Ninth IEEE International Conference on Data Mining*, pages 1–10. IEEE, 2009.
- [10] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [11] Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- [12] Young-Bum Kim, Dongchan Kim, Joo-Kyung Kim, and Ruhi Sarikaya. A scalable neural shortlisting-reranking approach for large-scale domain classification in natural language understanding. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers)*, pages 16–24, 2018.
- [13] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670. ACM, 2010.
- [14] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- [15] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214, 2011.
- [16] Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Online-to-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pages 1–9, 2012.
- [17] Alexandra Carpentier and Rémi Munos. Bandit theory meets compressed sensing for high dimensional stochastic linear bandit. In *Artificial Intelligence and Statistics*, pages 190–198, 2012.
- [18] Hamsa Bastani and Mohsen Bayati. Online decision-making with high-dimensional covariates. Available at SSRN 2661896, 2015.

- [19] Djallel Bouneffouf, Irina Rish, Guillermo A Cecchi, and Raphaël Féraud. Context attentive bandits: contextual bandit with restricted context. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 1468–1475. AAAI Press, 2017.
- [20] Lijing Qin, Shouyuan Chen, and Xiaoyan Zhu. Contextual combinatorial bandit and its application on diversified online recommendation. In *Proceedings of the 2014 SIAM International Conference on Data Mining*, pages 461–469. SIAM, 2014.
- [21] Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- [22] Ashkan Sharabiani, Adam Bress, Elnaz Douzali, and Houshang Darabi. Revisiting warfarin dosing using machine learning techniques. *Computational and mathematical methods in medicine*, 2015, 2015.
- [23] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.