

# CAL: A Framework for Managing Consistency, Availability, and Battery Lifetime in Mobile Applications

## ABSTRACT

We propose a framework that allows mobile applications to continuously choose between consistency, availability, and battery lifetime as network bandwidth and availability change and battery levels drain. Our system, CAL, provides various mechanisms and interfaces so applications can adjust their performance along these three axes in the tradeoff space as operating conditions change. We analyze a month's worth of mobile and WIFI network-coverage data and characterize the variability of the underlying network and discuss the implications on mobile application design. We describe how CAL eases application design in this context and demonstrate it in an indoor, energy-auditing application – an application for collecting information and querying the physical environment for live energy consumption feeds in buildings. We introduce a new prefetching technique, aimed to optimize fetch size and object set according physical proximity of the objects as inferred by object reference times. We show that CAL helps provide higher consistency and availability at lower energy cost, and degrades gracefully as resources become exhausted.

## 1. INTRODUCTION

A new class of mobile applications is emerging on mobile phones. The new frontier for enabling ubiquitous connectivity is moving indoors. Applications are becoming more reliant on ubiquitous, continuous connectivity, however there are fundamental challenges to enable an uninterrupted user experience in an indoor setting.

A variant of participatory sensing, whereby we're not just collecting information about the environment, but we expect the participants to be the authority in the construction of the environment.

This work will have similarities to Informed Mobile Prefetching (IMP, mobisys2012).

## 2. MOBILE CONNECTIVITY CHARACTERISTICS

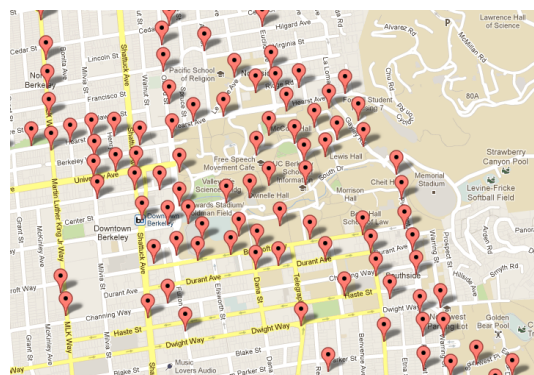


Figure 1: Locations.

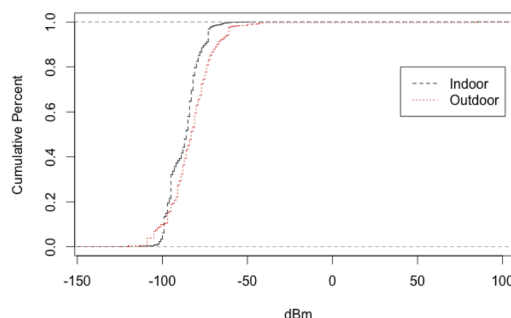
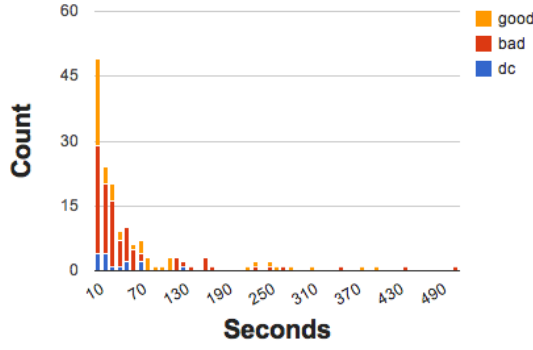


Figure 2: CDF for indoor and outdoor mobile signal strength measured indoors and outdoors. The distributions are very close and only slightly different in strength.

### 2.1 Methodology

We collected approximately 10-days worth of data traces from 6 participants in our lab using an application called the 'Network Signal Info Pro'. The application logs both mobile network and wifi network connectivity, passively, taking samples of the network state every 10 seconds. In order to deter-



**Figure 3: Percentage of time on each network.**

mine the connectivity characteristics indoors versus outdoors, we filtered the data points by location, where a location is defined by a center and a 100 meter boundary around that location. If the user remains at the location for long then 45 minutes, they classify that location to be an ‘indoor’ location, all other locations are ‘outdoor’ locations. We verified our indoor locations by inspecting them their latitude, longitude coordinates on a map.

For each class (indoor/outdoor), we examine the time-of-day connection characteristics. We specifically examine the distribution of signal strength at different times of day. We also look at the frequency and length of disconnection events.

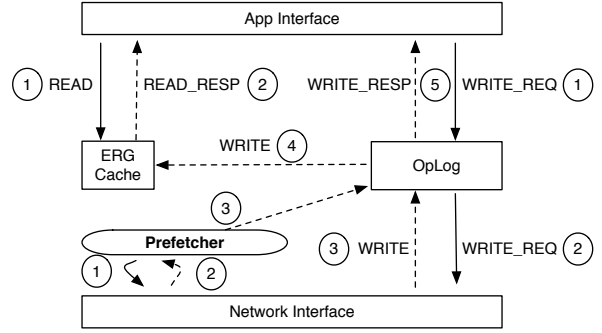
## 2.2 General connectivity and access

## 3. SYSTEM DESIGN

Our system is designed to provide a set of basic services and an API that application designers can use to both reason about and implement strategies for dealing with 1) consistency 2) availability 3) and battery lifetime. We provide an *simple API* that gives the application a way to make local decisions based on the age of application-level objects, a *Prefetcher* that prefetches data at an adjustable rate, and a *cache* that is used to deal with both reduced latency and increased availability. Much of the library consists of a set of interfaces, that encapsulate the notion of an application-level object, the application server, the operations that can be performed on those objects, and an *Expression*, which consists of a set of operations, executed in a given order, atomically.

To deal with disconnected operations, we provide a *OpLog*, which logs the set of operation that

were performed locally to application objects, that should eventually be performed on the server. Many of our API calls allow the user to submit a callback object, which is triggered when the certain operation is performed on the object on the server. This two-phase approach is necessary to keep the application responsive in times where access to the network is down. We also allow give the user hooks into the synchronization process, which monitors the consistency between the local object copies and those on the server.



**Figure 4: Standard mechanisms for consistency management on the phone. All READ request go to the local cached version of the ERG. All WRITES must go through the OpLog.**

Figure 4 shows a high-level overview of the components in the architecture and how they interact. We discuss these components in the following sections. We start with the library that runs on the client and the optional library to run on the server. We also discuss the benefits and implications of using either just the client or both the client and server stubs, particularly with respect to the consistency model.

## 4. CLIENT STUB

### 4.1 API design rationale

Different mobile applications have different consistency, energy, and responsiveness requirements. Unlike traditional application, mobile application need to be able to continuously adjust their approach to fall on different points in the trade-off space along each of these axes, throughout the lifetime of the application. Our goal is to design an API where the application programmer can reason about and implement applications that dynamically adjust their

API call	Description
read(objectName, save)	Reads the object from the server if a connection is available. Otherwise, reads from the local cache.
read(objectName, freshness, save)	Reads the object from local cache if it's $\leq$ freshness time units old. Otherwise, reads from server.
read(objectName, callback, save)	Reads from the server when the server is available. Callback is triggered after the fetch is complete.
readEOP(objectName, save)	Reads the object from the server if a connection is available and its within the energy budget. Otherwise, reads from the local cache.
readEOP(objectName, freshness, save)	Reads the object from local cache if it's $\leq$ freshness time units old. Otherwise, reads from server.
readEOP(objectName, callback, save)	Reads from the server when the server is available. Callback is triggered after the fetch is complete.
write(objectName, data, op)	Sends a request to the server to apply the operation remotely and copies the object locally after the operation complete. If the server is unavailable, applies the operation locally and logs it.
write(objectName, data, op, freshness)	Writes the operation to local copy if it is $\leq$ freshness time units. Otherwise tries to write to the server.
write(objectName, data, op, callback)	Writes to the server. Callback is triggered when the write is complete. A local copy is cached.
writeEOP(objectName, data, op)	Sends a request to the server to apply the operation remotely and copies the object locally after the operation complete. If the server is unavailable, applies the operation locally and logs it.
writeEOP(objectName, data, op, freshness)	Writes the operation to local copy if it is $\leq$ freshness time units. Otherwise tries to write to the server.
writeEOP(objectName, data, op, callback)	Writes to the server. Callback is triggered when the write is complete. A local copy is cached.

**Table 1: Summary of the main API calls of the Context Object Layer. Each call allows the designer to reason about and implement along different points in the tradeoff space between consistency, responsiveness/availability, and energy consumption.**

position along these three tradeoff axes as the run-time conditions change.

#### 4.1.1 Consistency

Eventually consistency is the consistency model that can be reasonably be implemented in our system. Variable bandwidth and disconnection, as well as variable energy consumption policies constraint the consistency model to best effort, eventually consistent. All clients eventually see exactly the same sequence of changes of the values for each object. We draw an important semantic distinction between the use of the just the client stub and the client-server combination. Using only the client stub gives speculative eventual consistency. Operations are performed, proactively, on local copies of the object. At some point in the future, those updates will be applied to the server. If the server cannot execute the given operations, the server object is copied on the server and the application is notified.

Operations are an important components of the approach we take. We force the application designer to define the set of operations that can be performed by the application server, so that they may also be applied on local copies of objects. For some applications, sets of operations must be performed atomically, or not at all. We define such a collection as an *Expression*. Expressions are only supported if the server stub is implemented, since only the server can execute operational sets atomically and only the server can resolve conflicts between clients.

#### 4.1.2 Availability and responsiveness

Network connection availability and quality change through space and time. Most applications need to make in-time decisions about which data and how much data they are going to fetch for the application. For indoor, interactive application, the network conditions can vary substantially in different parts of the indoor environment. Our api must support a caching layer for responsiveness and proactive prefetching for availability. We aim to allow the user to vary the size and frequency of the prefetch mechanism, thereby allowing for dynamic adjustment of application performance and availability as conditions dictate.

#### 4.1.3 Battery lifetime

We can attain high levels of consistency at a higher average cost of energy, by making our all writes go directly to the application server. For each read/write of an application object, there's an associated, time-varying, cost. Our API should allow applications to make the appropriate choice, given the

current operational conditions. For example, we can achieve high levels of consistency if all reads/writes are write-through to the server. At the other end of the spectrum, we do all reads/writes from locally. In the latter, we run the risk of dealing with stale objects.

Table 1 shows the basic API that is made available to the runtime of the client application. We essentially support three types of calls, each with three different sets of parameter and their own semantics.

## 5. FETCH COST

To determine the energy cost of a fetch, we need to calculate the time it takes to fetch the object.

We have variable connection quality and variable object size.

1. Maintain a weighted average of the size of the object to fetch.
2. Maintain a weighted average of the connection speed.
3. Maintain availability of each network based on the number distribution of object transfers over particular networks.

## 6. ENERGY BUDGETING

The energy budgeter keeps track of the energy consumed by this application using the same energy model formulated by Balasubramanian et al. [1].

## 7. APPLICATION EXAMPLE: ENERGY LENS

### 7.1 Access pattern and available

## 8. RELATED WORK

Bayou [4] is a system architecture that is designed to provide high-availability, and variable consistency while supporting application-specific conflict resolution and application-controlled inconsistency management. The system was specifically designed for mobile clients with variable connection quality and periodic disconnection. In many ways, the design goals of Bayou are very similar to ours. We both aim to provide highly available read and write operations, we both rely on an eventual consistency model, both support application-level detection of conflicts, application-specific conflict resolution, permitting disconnected clients to see their own updates. In contrast, we do not explicitly provide session guarantees. We assume all clients are

writing to the same application servers and that updates to the server by any client is consistency across a cluster of application servers.

Our system, however, introduces the notion of variable consistency and availability with respect to energy consumption. Typically applications are willing to spend more energy per I/O operation if connectivity is available and the cost/bit is cheap. Moreover, we include a prefetcher in our framework that adjusts its prefetch download size based on time-based hierarchical clustering of object reference time, dictated largely by the proximity of physical objects to the user interacting with the application. We focus on mobile applications that are interacting with objects in the physical world, not necessarily a virtual object on a server being edited by a random set of clients. Most reads on the virtual object will come from clients near that object in physical space, thereby limiting the number of concurrent interactions and the prefetch size to the space physical space occupied by those client’s primary users. Bayou is a database system for mobile clients, we are mainly client-side framework with optional server-side support to stronger consistency and ordering guarantees.

PocketWeb [3] is an prefetch framework for mobile phones that prefetching dynamic web content. PocketWeb uses a machine learning approach based on stochastic gradient boosting techniques to model the mobile web browsing patterns of mobile users. The main observation that allows their system to effectively prefetch content is that there are user-specific spatio-temporal access patterns. Their technique builds a model per user and prefetches 80%-90% of the content for 60% of users 2 minutes before they access it. Proactive, periodic prefetching is necessary dynamic web content, in particular. By limiting our framework to the class of indoor, interactive mobile application – where the user is directly interacting with the space around her through her phone – we limit the scope and access pattern to the spatial proximity to the world around them. With the increase in embedded sensing, we will see more applications that fall into this class, and prefetching as many virtual objects near the user as possible, will become more important. We also allow the application designer to control the frequency and scope of the prefetch. Fetching larger sets more frequently when connectivity is good and energy level is high.

Informed Mobile Prefetch (IMP) [2] using an adaptive strategy for prefetching user-specified application-level objects. IMP notes the objects that are referred to most often and treats the connection and

energy-level as a resource to be budgeted depending on the anticipated demand for those object and the projected energy and bandwidth supply. **We use the same prefetch decision criteria, except that we consider the cost benefit of fetching on a collection of items that considers the level in the hierarchical cluster tree, rather than calculate it on a per-item basis.**

## 9. TRACED-DRIVEN ANALYSIS

For our experiments, we ran a trace-driven analysis using traces from a mobile energy application conducted in a campus building. It includes the registration of new devices, meters, and spaces throughout the building as well as scans for information about those devices, from occupants in the building that want to learn more about the power-draw history of the devices. The audit was conducted by multiple students simultaneously, each walking around the different floor in the building, registering new devices and inputting information about each device. Because the traces do not include connectivity information, we take the indoor connectivity traces and simulate the effect of those conditions on the application traces.

We examine the consistency, availability and relative energy consumption of the application with and without the use of CAL.

We measure the following:

- Percentage of time the application is unavailable.
- The total amount of energy consumed by the
  - something
- Graph of object id and absolute access time (by offset from first access).
- Graph of the trace based on object id and color each cell according to the closeness in time for the weight absolute difference in object access times.

### 9.1 Energy Lens Application Trace-Based Simulation

We collected several months of logs from a mobile energy audit application summarized here:

1. How many unique user scans?
2. How many items were scanned in sequence per user?
3. How many moves?

4. How many removals?
5. How many updates?

We extracted the portion of the trace during which an audit took place. With the original version of the application an audit event must be aborted if the connection to the network becomes unavailable. We modified the application using CAL, allowing auditors to make local updates to the state of the audit and update the server when the connection becomes available. The portion of the audit trace that we use took place on a single floor. For energy audits in larger, older buildings in locations with poor coverage, we expect that battery consumption will become an issue. Moreover, auditors typically audit more than one building while on assignment [?].

We used these as the basis for a simulation of the application with multiple simultaneous users. We look to observe the effects of increased contention for physical interaction and virtual updates of the item observe how CAL helps manage that contention along the aforementioned tradeoff space. With connection-quality log traces to simulated the effects of varying connectivity.

1. Items scans in a room.
2. Item tag updates.
3. Item movement.
4. Item removal.
5. Probability of going to another room.
6. Probability of going to another floor.

Essentially we’re simulating a many, simultaneous of readers and writers.

## 9.2 Newsreader trace

To demonstrate how CAL performs for other application, we construct a newsfeed trace. The newsfeeds trace is constructed in a as described by Higgins et al. [2]. Google Reader provides statistics on the last 30 days of feeds with the average.

*Take the most popular feeds, find the median article size, assume each article is read from a uniform distribution between 30 and 60 seconds per article and then fetch the next one. We compare the response time, energy consumption, and availability. We simulate the loss of internet connectivity and energy level? Better consistency and availability (leverage the opportunity to prefetch an article.) What about energy? We’re not necessarily minimizing energy consumption, we’re managing the energy budget.*

*We want to demonstrate how the application can choose between the 3 axes by changing the calls in the API.*

*We also want to show that the prefetching strategy works well for a certain class of applications – the ones driven by physical-object-based scanning of the environment.*

## 10. REFERENCES

- [1] N. Balasubramanian, A. Balasubramanian, and A. Venkataramani. Energy consumption in mobile phones: a measurement study and implications for network applications. In *Proceedings of the 9th ACM SIGCOMM conference on Internet measurement conference*, IMC ’09, pages 280–293, New York, NY, USA, 2009. ACM.
- [2] B. D. Higgins, J. Flinn, T. J. Giuli, B. Noble, C. Peplin, and D. Watson. Informed mobile prefetching. In *Proceedings of the 10th international conference on Mobile systems, applications, and services*, MobiSys ’12, pages 155–168, New York, NY, USA, 2012. ACM.
- [3] D. Lymberopoulos, O. Riva, K. Strauss, A. Mittal, and A. Ntoulas. Pocketweb: instant web browsing for mobile devices. In *Proceedings of the seventeenth international conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS ’12, pages 1–12, New York, NY, USA, 2012. ACM.
- [4] D. Terry, M. Theimer, K. Petersen, A. Demers, M. Spreitzer, and C. Hauser. Managing update conflicts in bayou, a weakly connected replicated storage system. In *Proceedings of the Fifteenth ACM Symposium on Operating Systems Principles*, pages 172–183, 1995.