

Anders Emil Nielsen



Kongens Lyngby 2016

Technical University of Denmark
Department of Applied Mathematics and Computer Science
Richard Petersens Plads, building 324,
2800 Kongens Lyngby, Denmark
Phone +45 4525 3031
compute@compute.dtu.dk
www.compute.dtu.dk

Summary (English)

The goal of the thesis is to ...

Summary (Danish)

Målet for denne afhandling er at ...

Preface

This thesis was prepared at DTU Compute in fulfilment of the requirements for acquiring an M.Sc. in Engineering.

The thesis deals with ...

The thesis consists of ...

Lyngby, 04-July-2016



Not Real

Anders Emil Nielsen

Acknowledgements

I would like to thank my....

Contents

Summary (English)	i
Summary (Danish)	iii
Preface	v
Acknowledgements	vii
1 Introduction	1
1.1 Problem	2
1.2 Requirements	3
1.3 Context / Running Example	4
1.4 Contributions	5
1.5 Outline	5
2 Caching Model	7
2.1 Caching Basics	7
2.1.1 Architecture	9
2.1.2 Timeline Model	10
2.2 Evaluating Caching Techniques	11
3 Caching Approaches	15
3.1 Invalidation Techniques	15
3.1.1 Expiration-based Invalidation	16
3.1.2 Key-based invalidation	17
3.1.3 Trigger-based Invalidation	18
3.1.4 Trigger-based Invalidation with Asynchronous Update . .	22
3.1.5 Write-Through Invalidation	23
3.1.6 Automatic Invalidation	24

3.1.7	HTTP Caching	29
3.2	Choosing the Right Caching Technique	29
4	Smache: Cachable Functions	31
4.1	Existing Approaches are Not Sufficient	31
4.2	The Cachable Function Model	34
4.2.1	Restricted to Pure Functions	34
4.2.2	Making Functions Cachable	34
4.2.3	Cache Object Localization	35
4.2.4	Automatic Cache Invalidation	35
4.2.5	Data Update Propagation	36
4.3	Implementing Cachable Functions in Python	36
4.4	Discussion	38
5	Automatic Cache Invalidation	39
5.1	Simple Object Dependence Graph	40
5.2	Dependency Data Structure for Cachable Functions	41
5.2.1	Declaration Dependence Graph	41
5.2.2	Instance Dependence Graph	44
5.3	Dependency Registration	45
5.4	Invalidation Propagation	46
5.4.1	Timestamp Invalidation	47
5.4.2	Database Wrapper Triggers	49
5.5	Implementing Automatic Invalidation	49
5.6	Discussion	49
6	Data Update Propagation	51
6.1	Existing Data Update Propagation Approaches	51
6.2	Race Condition on Write-Through Invalidation	52
6.3	The Data Update Propagation Algorithm	52
6.4	Implementing the Data Update Propagation Algorithm	54
7	Results and Evaluation	57
8	Conclusion	59
8.1	Future Work	59
A	Stuff	61
A.1	Code Snippet for Trigger-based Invalidation with Asynchronous Update	62
	Bibliography	65

CHAPTER 1

Introduction

“There are only two hard things in Computer Science: cache invalidation and naming things.”

– Phil Karlton

Web applications are becoming more and more dynamic with more personalized content that often requires complex data queries or computations based on large amounts of data. These computations can become a performance bottleneck in the application, which leads to slow response times and poor user experience for the users.

The performance can often be optimized by profiling and analyzing the code behind the computation, but this is often not the easiest solution and in some cases the complexity or amount of data used makes it difficult to achieve a satisfactory performance. Caching is a popular solution for improving the performance and scalability in these cases since it allows for a simple, scalable and generic way of addressing bottlenecks in the web applications.

Although it sounds like a silver bullet it also places a burden on the programmer that must locate and update the cached values while preserving consistency guarantees. This challenge is for example seen in an outage of the whole Facebook system:

The intent of the automated system is to check for configuration values that are invalid in the cache and replace them with updated values from the persistent store. This works well for a transient problem with the cache, but it doesn't work when the persistent store is invalid.

Today we made a change to the persistent copy of a configuration value that was interpreted as invalid. This meant that every single client saw the invalid value and attempted to fix it. Because the fix involves making a query to a cluster of databases, that cluster was quickly overwhelmed by hundreds of thousands of queries a second.

Robert Johnson [Joh10]

This example shows how critical the caching system can be and the importance of correctness.

This thesis will address this issue by researching the latest caching technique proposed in research and used in practice and contribute with a design and implementation of a caching system in the Python programming language.

1.1 Problem

Most of the existing caching solutions are based on a pull based caching strategy, where the computation runs and the cached value is stored when the client requests the cached value. After the result has been computed and cached, the client will be presented with the cached value until the it is invalidated.

The pull based caching strategy has the advantage that we only have to cache content that is being used, but it also means that the first time a client asks for the value, it has to wait for the computation to finish. This is not optimal with relation to user experience since the user has to wait in order to be presented with the requested content. To solve this problem we have to precompute the cached values such that the user is presented with content as soon as it is requested.

Besides the performance problem on the initial request, existing caching solutions leaves responsibility for the programmers to maintain the cache in order ensure consistency and an appropriate level of freshness.

These problems with existing caching solutions, presents two major challenges:

Cache Management

The first challenge related to cache management is faced in any caching system, where the programmer has to manage the caching system by naming the cached value and keeping them up to date such that the users are not presented with unexpected content.

One particular challenge within cache management is *cache invalidation* that relies on the programmer correctly identifying every underlying data that affects the given cached value. The programmer then has to declare a way for the cached value to be invalidated when any of the underlying data changes. This analysis is difficult since it requires global reasoning about how the underlying data changes in the web application and which computations are cached. Furthermore if the computation behind the cached value is altered to depend on new underlying data, the cache invalidation also has to change, making the cache prone to errors if the latter is forgotten.

We discuss this more in chapter 3 and 5.

Data Update Propagation

The second challenge is related to the task of efficiently keeping the cached values fresh while ensuring the consistency between the cache and the storage system. This challenge will be addressed in chapter 6.

1.2 Requirements

The final solution addressing the problems described, will be designed with the following non-functional requirements:

Software design: Must be designed to be maintainable such that the programmer that uses the caching system understands how it works from using it and has the ability to extend it. The design of the system should also be flexible to support multiple storage systems and caches.

Adaptability: Should be convenient and easy to adapt into existing systems. Furthermore the usage of the system should be easy to understand.

Efficiency: Should be efficient with relation to performance such that it does not make existing operations of the systems significantly slower. It should also

be efficient with relation to the system load such that it does not use more computational power than necessary to achieve the goal of the system.

Scalability: Should be designed for scalability in the sense that the design should still be efficient for large amount of data and correct when the web application is scaled horizontally.

Fault-Tolerance: Should be designed with considerations on reliability, availability, integrity and maintainability.

1.3 Context / Running Example

The problem and requirements are based on a running web applications - the Peergrade.io-platform to ensure that the system is also designed, implemented and tested to be used in practice.

Peergrade.io is a platform for facilitating peer-evaluation in university and high school courses. Currently the platforms serves multiple institutions and thousands of students. One of the key parts of the platform is showing various statistics about the performance of the students in a course. The statistics are based on advanced calculations which take up a large amount of time and needs to be recalculated on small changes to the underlying data.

To relate the solution to a practical example, the thesis will use the following code as a running example:

```
1 def course_score(course)
2     participants = Database.find_all_participants_in_course(course)
3     total_score = 0
4     for participant in participants:
5         total_score += time_consuming_participant_score(participant)
6     return total_score / len(participants)
7
8 def time_consuming_participant_score(participant):
9     grades = Database.find_all_grades_for_participant(participant)
10    return numpy.advanced_statistical_method(participant)
```

Figure 1.1: Code with the running example written in Python

In this example we have a function `course_score` that computes the average score in a given course by fetching the participants from the primary data store

and through iterations of each participant calculate the average score using the function `participant_score` that calculates the score of a single participant using a long running external method.

1.4 Contributions

This thesis addresses challenges described in section 1.1 and requirements in section 1.2 in the context of a caching system. The result will be a design of a cache system solving those challenges based on the requirements. The design will be implemented in Python and made available as open source to allow further research and extensions from the implementation. The implementation will therefore also have a focus on delivering a maintainable and tested library.

1.5 Outline

This thesis is structured as following:

With the motivation, problem, and requirements described in this introductory chapter, chapter 2 will give an introduction to the basics of caching and present the models used and a set of criteria used to evaluate caching approaches.

The following chapter 3 will describe the existing caching approaches described in literature and used in practice.

Based on the knowledge of existing solution, chapter 4, 5, and 6 describes the solution suggested by this thesis. The solution is explained using the programming model in chapter 4, which is written as a self-contained part of the solution. This model is then extended with automatic invalidation in chapter 5 and afterwards with in-place updates in chapter 6.

The solution is then evaluated and discussed based on test results and the initial requirements in chapter 7.

Chapter 2 introduces caching by presenting the basic caching algorithm, the common architecture of caching system, the models used to describe caching approaches, and criteria used to evaluate caching approaches to find the appropriate technique.

Chapter 3 will describe existing caching approaches with relation to caching evaluation criteria and explain how to choose a caching technique based on these criteria.

Chapter 4 presents the solution suggested by this thesis for the given context and requirements. First the techniques of existing approaches are discussed and followed by a description of the programming model - cachable functions.

Chapter 5 describes in more details how the cachable functions are extended to have automatic invalidations using declared dependencies to underlying data.

Chapter 6 explains the data update propagation algorithm used to extend the cachable functions to have in-place cache updates.

Chapter 7 goes through the test results related to performance, efficiency, and cache memory usages. The solution are then discussed and evaluated based on the test results and the initial requirements.

Chapter 8 finalize the thesis with a conclusion.

CHAPTER 2

Caching Model

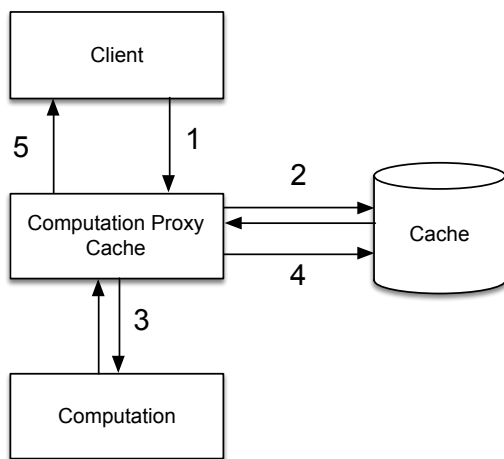
In order to get a common understanding of caching and the terminology related to the topic, we will go through the basics of caching by describing the architecture of a web system using a cache, present a model based on a timeline that introduces the different events involved in caching and last we list criteria for evaluating a caching technique in order to choose an appropriate technique for a given use case.

2.1 Caching Basics

In general caching is about storing the result of a computation at a where you are able to retrieve it fast, such that it is possible to get the result fast instead of recomputing it. This basic algorithm is illustrated on figure 2.1 and can be described as following:

In some cases, where the client is not allowed to wait for the computation to run, step 3-4 are replaced by a step that simply returns an empty value.

If we look at the cached object from an abstract point of view, we can see it as a *result of a function* given certain *inputs*. Sometimes the inputs are data



1. The result v of a computation f is requested.
2. If v is cached and is valid, we go to 5 with $v' = v$. Else we continue.
3. We run the computation f
4. The new result v' of f is stored in the cache.
5. v' is returned.

Figure 2.1: The flow of basic caching

from a storage system, sometimes it's the result of an API call to some external resource, sometimes it's global variables in the code. These *inputs* we from now on be referred to as *underlying data*.

In order for the algorithm to work, we need to be sure that when we store the cached object it has to be uniquely identified by some key such that when we lookup the value as in step 2 of the algorithm, we always locate v . This presents one of the challenges of cache management, which is in many cases closely related to cache invalidation (one of the two hard things in computer science¹).

In the algorithm cache invalidation is simply described as a *check*, but in reality this is the hard challenge of caching. The *check* could be a precomputed indicator from earlier triggers, it could be based on a key derived from the underlying

¹Not scientifically, but at least a favorite saying of Martin Fowler and quote by Phil Karlton

data or some timestamp. These cache invalidation approaches will be described more in section 3.1.

2.1.1 Architecture

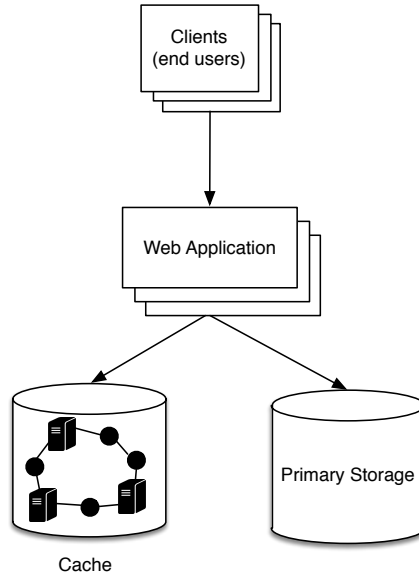


Figure 2.2: The assumed architecture of the system

The architecture of the web application in which the cache is used is important to how the cache system. We assume that the architecture of the system is a common web application architecture as illustrated in figure 2.2 consisting of a web application that serves HTTP-requests from the client (represented by the user) and interacts with a primary storage database to store and load data. To store and fetch the cached content we introduce a cache database. This could easily be the same unit as the primary storage database, but we separate them for better clarification.

Most modern web applications need to serve multiple users at the same time, which means the web application must run on multiple processes² either on a single or multiple machines. We will therefore treat the web application as a distributed system.

²This is processes as an abstract term used in distributed systems. If we need to be implementation specific this could just as well be threads.

A popular choice for cache databases are key-value stores such as Redis ³ and Memcached ⁴ since they are simple distributed key-value stores that lives in memory and therefore allows for scalability and high-performance operations. In order to support most practical web applications, we will therefore make the assumption that the cache database has the same functionality: it should be possible to store arbitrary content for a given key. Furthermore the final solution of this thesis require the cache database to support atomic transactions for a given key, which are supported in Memcached using the CAS command [Mula] and in Redis using the WATCH-command [Mulb]. The reasons behind this assumption will be explained more in chapter 6.

2.1.2 Timeline Model

As with the algorithm on figure 2.1, caching can be described by a series of events. The ordering of events decides whether the cached content or a fresh computation is presented to the client. To describe the different caching techniques, we will use a timeline model with a stream of events. One timeline describes the events occurring in a single process. We can therefore assume that there exists a total ordering of events for a single timeline. To be able to describe the caching techniques explained in this thesis, we will define the following events:

- **Requested** is the event occurring when a client requests a given cached object
- **Computation Started** is when the computations is started
- **Computation Finished** is when the computation is finished and a result is returned
- **Stored** is happening when a given cached object is stored in the cache database
- **UD Updated** is when the underlying data has been updated
- **Invalidated** is when the system has knowledge that the cached object is invalidated

Alongside these events the timeline model illustrates the interval in which the given cached object is considered valid (the cache system will serve the cached

³<http://redis.io/>

⁴<https://memcached.org/>

value) and when it is actually fresh (consistent with underlying data). The validity interval illustrates when the cache system will respond immediately, which means if the cached object is always considered valid then the approach supports immediate responses. In the time interval, where the fresh and valid interval overlaps, the approach will serve a fresh version of the cached value. If they always overlap then the given approach will guarantee strict freshness.

The timeline model applied to the basic caching algorithm 2.1 is illustrated on figure 2.3.

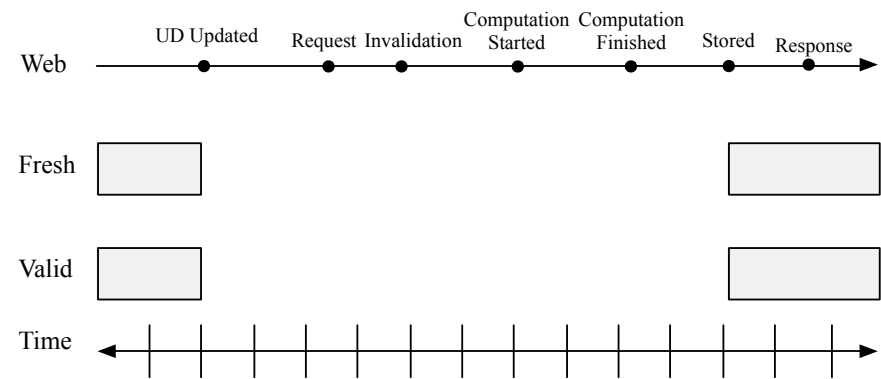


Figure 2.3: The timeline model applied to the basic caching algorithm

2.2 Evaluating Caching Techniques

To choose the correct caching technique for a given use case, we need to know the criteria for finally picking the best suited. The overall goal of caching in web development is to achieve a better user experience by getting a better performance and to save money by using less CPU power. If we assume that the cache system is able to retrieve cached objects fast, the goals can be achieved by “hitting the cache” as often as possible. We can measure this using the metric *cache hit rate*, which refers to the rate at which the cache is hit among the total number of requests for the cached object.

To evaluate the caching technique for a given use case, we will see the situation from the perspective of the client (i.e. a user). The client makes a request for some content that is served by the web server. This content contains of one or more computations that can be cached individually.

We will not make any assumptions about the content send by the server, which means the content could be the result of multiple cached computations. Each of these computations are based on some underlying data e.g. served from the primary store. In the case where some computations are based on the same data we have to keep the different results consistent. Furthermore the response might be based on both cached content and content loaded directly from the primary store in which case we have to keep the data consistent across the cache database and the primary storage. This issue leads to the criteria of the level of consistency.

There exists multiple levels of consistency, but to keep it relevant and simple, we will evaluate the level of consistency using a binary value: either the caching technique ensures consistency with the data from the primary storage or else it doesn't.

Another parameter is the freshness of the content returned by the cache. It is most desirable to have content that is as fresh as possible as oppose to having stale data, but in the end the goal is to make the served content make sense for the user. Consider example 2.1 of the Peergrade.io platform, where there are two types of users. If the teacher changed the description of an assignment and a student requested and read the description 1 min after then it would not be unexpected behaviour to show the old description from the students point of view. On the other hand if the teacher requested the description 1 min after, it would be unexpected behaviour to show the teacher an old version, since the teacher would think the description hasn't been updated. We will represent freshness as a binary parameter that we call "Strict Freshness", which evaluates whether a given cached object that is fetched is guaranteed to be fresh.

EXAMPLE 2.1 *At Peergrade.io the web application has a teacher and a student interface. Teachers create assignments and sees statistics about the grades students have given to each other through their feedback. The students see the assignments created by their teachers, are able to hand in their assignments, and grade other students' hand in.*

While the freshness describes what is expected behaviour with relation to the content, there also exists time limits with relation to keeping the user focused on the task. Miller and Card et. al. [Mil68, CRM91] describes these limits as:

- When the response time is **0.1 second** the user feels that the system is **reacting instantaneously**.
- A response time above **1 second** will **interrupt the user's flow of thought**.

- **10 seconds** is limit related to **keeping the user's attention** on the given dialog.

In the basic caching algorithm described on figure 2.1, the computation has to run after it has been invalidated. If the caching technique runs this computation in the same process as the request from the client, the client has to wait for the computation to finish in order to show the content to the user. If the time taken to compute the value exceeds the accepted response time limits described above, it could become critical for the user experience. Based on this fact we will introduce the binary parameter of whether or not the client has to wait for the computation to finish after invalidation. The weight of this parameter is proportional to the cache miss rate since it's only requests, which results in a cache miss that are affected by the slow response time.

It is not possible to both serve the cached objects immediately and have strict freshness. This can be proved using the example where a cached object is invalidated just before it is requested. We can only start computing the value at the moment it has been invalidated and given that the time between the invalidation and the request is smaller than the time taken to compute the value, it is not possible.

In cases where we choose immediate response time and we can tolerate serving cached objects that are not strictly fresh, we might want to keep the cached objects as up to date as possible. This means that we start computing the cached objects as soon as they have been invalidated. In order to achieve this we need some mechanism that invalidates depending cached objects when underlying data are updated. We will also represent "Automatic Invalidation" as a binary parameter.

So far we've considered parameters from a user's perspective, but to evaluate the caching technique fully we also need to see it from the perspective of the programmer since a big part of cache management is manually tracking dependencies between the cached content and the underlying data. We want the caching system to be transparent such that it is easy to add and remove caching from existing computations. Additionally we want the caching system to be robust such that when new code involving a cached computation is added, it should behave as expected without introducing errors. We cannot measure the level of robustness so we will therefore introduce the binary criteria: does the caching technique involve maintenance with relation to invalidation (shortened as invalidation maintenance)?

The last parameter we will introduce is also a requirement of the system: *adaptability*. Since adaptability cannot be evaluated objectively we will introduce a

scale of *high*, *medium* and *low*, where high means that is adaptable to most systems and low means that it is adaptable to a few systems. We will define the scale as following:

- **Low:** The technique has assumptions about the primary data store, cache database and/or application that means it cannot be applied to other common technologies.
- **Medium:** The technique is advanced and requires a lot of implementation effort or external libraries/processes to work.
- **Low:** The technique can easily be implemented and applied to existing applications.

The reason behind the *Low* criteria is that when the technique has assumptions about the technology behind the application it means the system is tied to using specific technologies, which makes it difficult to change when the given technology is obsolete or the requirements of the system change. If this is not considered a problem *Low* and *Medium* can be considered the same level of adaptability.

From this discussion, we can sum up the evaluation parameters as follows:

- **Consistency:** The cached object must be consistent with the other data
- **Strict Freshness:** The cached object must be as fresh as the state of the primary store when the value was requested
- **Automatic Update:** The cached object is automatically updated after it has been invalidated
- **Always Immediate Response:** The cached object must be served immediately after it has been requested
- **Invalidation Management:** Whether the programmer has the responsibility of maintaining the invalidation of the cached objects
- **Adaptability:** How adaptable the caching technique is to existing systems

CHAPTER 3

Caching Approaches

Since caching is an approach widely used in practice there already exists multiple caching approaches described in literature, articles on the internet and in open source code. To discover existing solutions for the problems and challenges faced in the thesis, we will research and evaluate existing caching approaches. The approaches will be analyzed and evaluated based on the criteria described in section 2.2 and requirements from section 1.2.

As already described the most difficult part of caching is the invalidation. We will therefore start by describing existing invalidation techniques followed by how the cached values are updated and related problems. We will also describe some specific caching approaches related to web development and finish the chapter by giving an approach to picking an appropriate caching approach for a use case based on the caching criteria.

3.1 Invalidation Techniques

When the underlying data for a cached value is updated we consider the cached value as invalid. Although this sounds as a trivial part, the mechanism for invalidating the cached value can become complex depending on the requirements for

the use case. To be able to navigate the existing solutions we will cover the general invalidation techniques as well as the more advanced techniques described in literature.

3.1.1 Expiration-based Invalidation

In cases where it is accepted to keep the cached values stale up to a given time interval, we can use the expiration-based invalidation technique, which is the simplest invalidation mechanism since the invalidation only depends on time and not updates of underlying data. Using the expiration-based invalidation we give up consistency and the level of freshness depends on the time interval set by the programmer, but we are able to achieve immediate responses.

The expiration-based invalidation works by assigning a TTL (Time to Live) to the cached value. At some point when the TTL has expired, the cached value is invalidated. The responsibility of invalidation cached values with TTL is often placed on the cached database by piggybacking the TTL to the cached value when it is stored. This is for example supported by Redis and Memcached.

To give a sense of how the relation between when a cached value is considered valid (it is served from the cache) and when it is actually fresh, the timeline model has been applied to the expiration-based invalidation on figure 3.1.

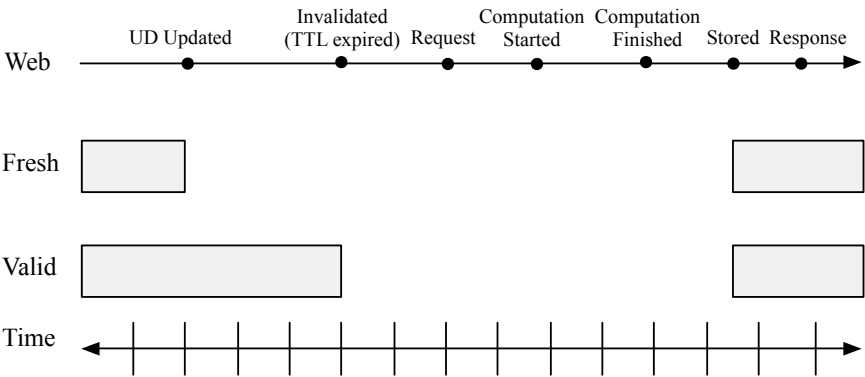


Figure 3.1: The lifecycle of the expiration-based invalidation technique

3.1.2 Key-based invalidation

When the cached values are required to be up to date with the primary storage, we must ensure that when underlying data is updated, the depending cached values must be recomputed before the cached value can be served. The key-based invalidation gives these guarantees by giving up immediate responses since the users have to wait for computations to finish if they request an invalidated cached value. As mentioned in section 2.2 the impact of this trade-off is proportional to the cache miss rate. This means that the key-based invalidation will not be suitable in cases where the computation time is too long or in cases where the cached values are updated too frequently.

Key-based invalidation works by constructing the cache key from parts of the underlying data such that the when the cached object should change, then the key also changes. [Han12]. The cached content is considered immutable and only have to be written once. This simplifies version management from the perspective of cache storage since there is no chance you read stale values if the key is assumed to be derived from the most recent version of underlying data.

The challenge of this method is to construct the key. To use this technique correctly (i.e. obtain the guarantees promised) the programmer must construct a key that is ensured to change when the cached value is considered stale. Furthermore to obtain a maximum hit rate, the key must not change when the cached value is fresh. Given that the key is optimally constructed, the timeline looks as on figure 3.2.

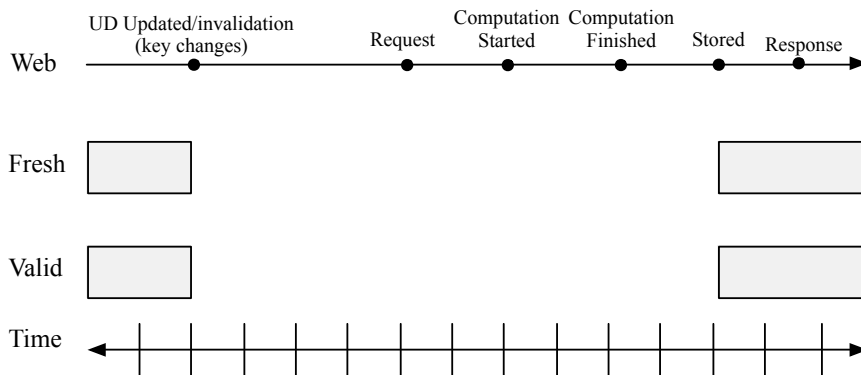


Figure 3.2: The lifecycle of the key-based invalidation technique

In the web application framework, Ruby on Rails, the key construction is simplified by using a key that includes the timestamps of the last update on some

underlying data. An example of this technique can be seen in code example 3.3. In this example we use the type, id and update timestamp as parts of the key. This means that the cached value is considered invalid if it is called with an entity that has a different type, id or update timestamp. The intuition behind these components is that we want a unique cache value for each entity and we want the value to be recomputed when the entity is updated.

The key-based invalidation performs invalidation at the moment, the cached value is requested, but it is considered invalid at the moment after the cache key components are updated (e.g. the *updated_timestamp* in the code example). This means we get interval the value is considered value and fresh are always overlapping, which means we have consistency.

A caveat of this method method is that it generates cache garbage since old versions of a cached value are not removed. To avoid the complexity of keeping track of the relations between the different, the responsibility for cleaning up is moved to the cache database. Fortunately cache databases (such as Redis and Memcached) implements different such cleanup algorithms that detects obsolete values based on some policy. One such policy called *Least Recently Used (LRU)* removes the cached values that are least recently used by keeping a timestamp of when the cached values where accessed last. Another policy called *Least Recently Used (LRU)* keeps track of the frequency in which the different values are accessed and removes the ones that are least frequently accessed.

3.1.3 Trigger-based Invalidation

Instead of invalidating the cached value when requested, the cached values can be invalidated based on certain events triggered when the underlying data is updated. Given that invalidation triggers are located at all the places where the underlying data is updated, we can achieve the same guarantees of consistency and freshness as with key-based invalidation. The timeline model of the trigger-based invalidation on figure 3.4 therefore looks similar to the key-based invalidation.

Where the key is used as an invalidation mechanism in key-based invalidation, the key is static for trigger-based invalidation i.e. the key for identifying a cached value does not change in its lifetime. The actual key becomes simpler since it is only responsible for uniquely identifying the cached value, but the localization is still a challenge since the key needs to be shared between the triggers and the places where the cached value is accessed. The responsibility of invalidating the keys are then moved to the definition of event triggers.

```
1 def time_consuming_participant_score(participant):
2     grades = Database.find_all_grades_for_participant(participant)
3     return numpy.advanced_statistical_method(grades)
4
5 def cache_key_for_participant_score(participant):
6     cache_key_components = [
7         'cached_participant_score',
8         participant.type,
9         participant.id,
10        participant.update_timestamp
11    ]
12    return '/'.join(cache_key_components)
13
14 def cached_participant_score(participant):
15
16     cache_key = cache_key_for_participant_score(participant)
17     if is_fresh_in_cache(cache_key):
18         return fetch_from_cache(cache_key)
19     else:
20         result = time_consuming_participant_score(participant)
21         set_cached_value(cache_key, result)
22         return result
23
24 # Load the participant from the primary storage
25 participant = ParticipantDB.load_one_from_database
26
27 # Call the cached version of the time_consuming_participant_score
28 print cached_participant_score(participant)
29
30 # This second time it is called, the return value for
31 # the time_consuming_participant_score cached, which
32 # means it returns the value from the cache
33 print cached_participant_score(participant)
```

Figure 3.3: Example of the key-based invalidation technique

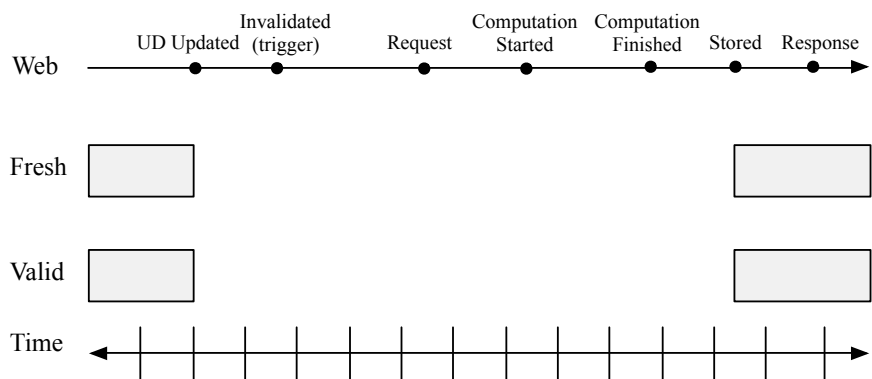


Figure 3.4: The lifecycle of the trigger-based invalidation technique

Having a static key also has the advantage that we can locate the old/stale value after invalidation as oppose to key-based invalidation where there is no relation between the versions of cached values. Using this information we can extend the technique to be more fault-tolerant by serving a stale value in the case where a computation fails¹. We can also extend the technique to be more flexible with relation to the properties as explained in section 3.1.4 and section 3.1.5.

The simplest type of triggers are manually defined code that invalidates a given key. A code snippet for a naive implementation of manual triggers can be seen in snippet 3.5. In practice the manual code triggers are often placed right after updates to the underlying data. Although this method is simple it often requires a lot of effort from the programmer and is prone to errors since it requires global reasoning of the application to identify the places where underlying data is updated.

Having a static key also introduces a challenge related to concurrency since we assume that there are multiple application servers. The challenge originates from the problem illustrated on figure 3.6, where an update to the underlying happens during the computation of a cached value. When the computation has finished it will incorrectly mark the cached value as fresh even though it is based on an old version of the underlying data, which makes it stale.

¹Here we assume that the application provides more value to the client by serving a stale value compared to serving an error or nothing


```
1 def time_consuming_participant_score(participant):
2     grades = Database.find_all_grades_for_participant(participant)
3     return numpy.advanced_statistical_method(grades)
4
5 def cache_key_for_participant_score(participant):
6     cache_key_components = [
7         'cached_participant_score',
8         participant.type,
9         participant.id
10    ]
11    return '/'.join(cache_key_components)
12
13 def cached_time_consuming_function(participant):
14     cache_key = cache_key_for_participant_score(participant)
15
16     if is_fresh_in_cache(cache_key):
17         return fetch_from_cache(cache_key)
18     else:
19         result = time_consuming_participant_score(participant)
20         set_cached_value(cache_key, result)
21         return result
22
23 # Load the participant from the primary storage
24 participant = ParticipantDB.load_one_from_database
25
26 # Call the cached version of the time_consuming_participant_score
27 print cached_participant_score(participant)
28
29 # This second time it is called, the return value for
30 # the time_consuming_participant_score is cached, which means it
31 # returns the value from the cache
32 print cached_participant_score(participant)
33
34 # Now we invalidate the cached value
35 cache_key = cache_key_for_participant_score(participant)
36 invalidate_cached_value(cache_key)
37
38 # Since the value has been invalidated
39 # the cached_time_consuming_function will recalculate
40 # the result when called at this point
41 print cached_participant_score(participant)
```

Figure 3.5: Example of how trigger based invalidation works with manual code invalidation

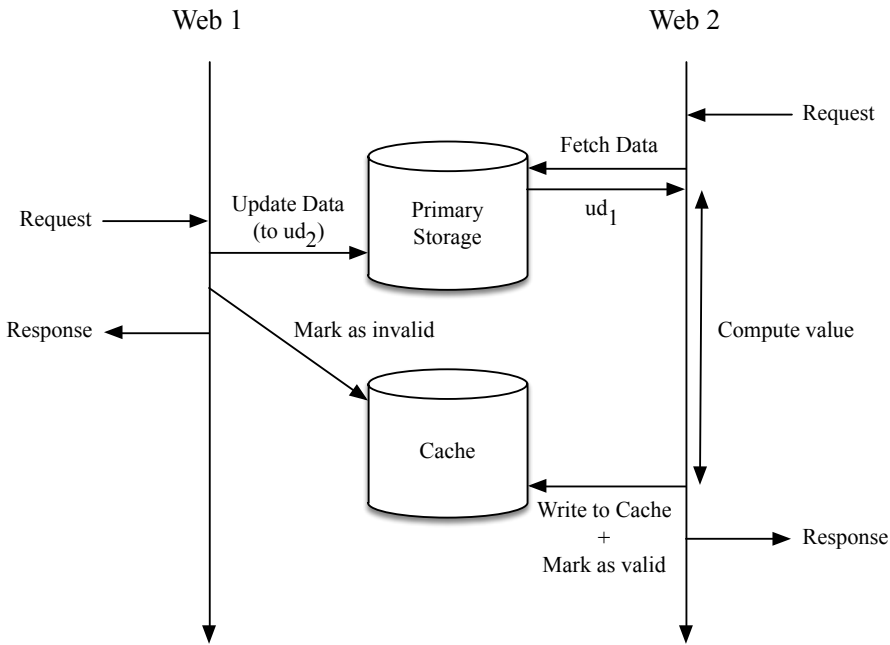


Figure 3.6: A scenario of the trigger-based invalidation that results in a race condition, where the cached value are being incorrectly marked as valid even though it is storing a stale value.

3.1.4 Trigger-based Invalidation with Asynchronous Update

We can extend the trigger-based invalidation by always serving the newest value from the cache and afterwards update the value asynchronously if it is stale. This way we always get an immediate response by giving up strict freshness and consistency. To give this guarantee fully the user have to wait for the computation the first time the value is requested if the application haven't "pre-heated" the cache i.e. included a build step before deployment that computes all cached values.

The timeline of this extension is as on figure 3.7. A naive implementation of this technique would be as the trigger-based seen in code snippet 3.5 with the modification that the system updates the value asynchronously instead of synchronously if no fresh value is found in the cache. The snippet for this can be found in appendix A.1.

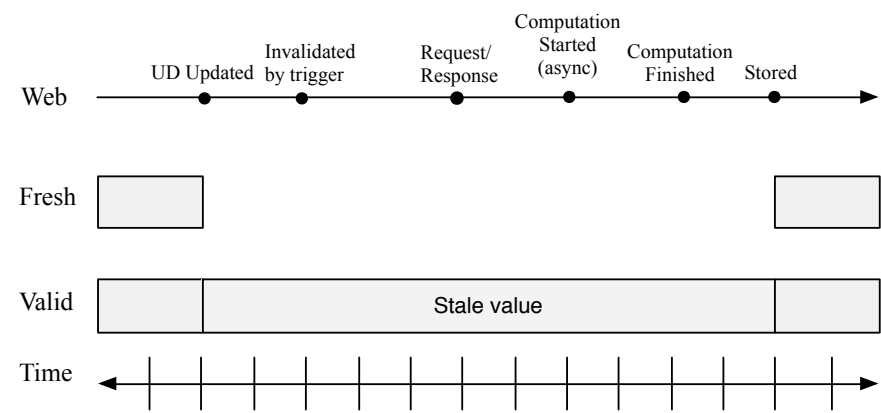


Figure 3.7: The lifecycle of the trigger-based invalidation technique where the value is updated in the asynchronous

3.1.5 Write-Through Invalidation

Write-through invalidation is also an extension to the trigger-based invalidation method, but instead of updating the cached values when the value is requested, the value is updated in the moment after it has been invalidated. This way we invalidate by writing through the existing version in the cache. The timeline model of this technique seen on figure 3.8 is similar to the asynchronous update extension, but the time interval in which it serves a stale value is only as long as the time taken to compute the value.

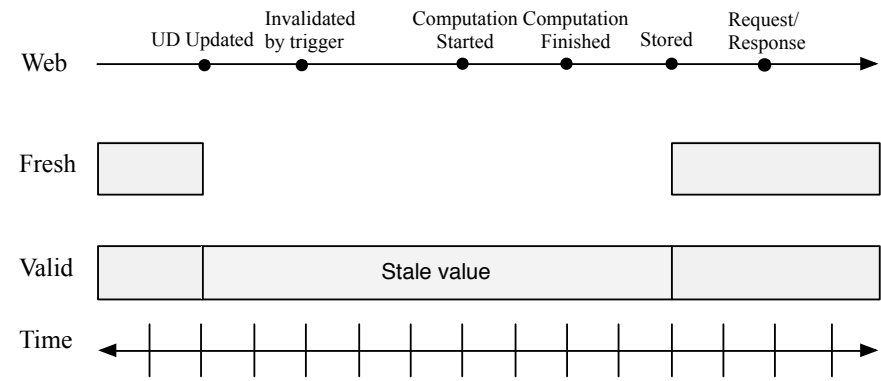


Figure 3.8: The lifecycle of the write-through invalidation technique

3.1.6 Automatic Invalidation

The trigger-based invalidation is an attractive technique since it can provide guarantees with relation to freshness and consistency while allowing for flexibility by extending it. But in practice the overhead for the programmer of managing the triggers and keeping integrity becomes a burden that makes it hard to maintain. A lot of research have therefore been done in making it easier to use trigger-based invalidation.

The cached objects are based on underlying data from the primary storage system, which means that changes to the underlying data also origins from the primary store. A lot of work have been put into using the database as the source of the triggers that invalidates the cached objects.

In [GZM11] the database wrapper ² is used to detect and trigger changes to the underlying data. This paper suggests a caching approach for caching database queries by declaring predefined queries with dependencies to underlying data using an extension to the database wrapper. The database wrapper is then responsible for detecting changes to underlying data and invalidate the affected cached queries. The advantage of using a database wrapper is that the caching system will still work if the database used by the wrapper is changed. On the other hand changes made to the database that are not made through the given database wrapper are not detected as a trigger, which leaves the responsibility of ensuring all changes are made through the database wrapper.

This problem can be solved using a database technology that is able to notify about changes directly from the database as done in the approach from IBM [CDI98, CID99] that uses the IBM DB2 ³ database. In this approach the triggers are intercepted by a cache-manager that is able to invalidate the affected values using a dependency graph. Having a cache-manager introduces a single point of failure and potentially a bottleneck, but it allows capturing dynamic dependencies, which is required for the given application. The details about this approach will be described in section 5.1.

In [Por12] a patch is made to the database such that the database is able to support so called “invalidation streams” that can help the caching system invalidate the affected cached values. The cache nodes stores information about how the different cache objects should be invalidated represented by invalidation tags. The invalidation stream from the storage system are then distributed across all cache nodes and used to invalidate the affected cached objects using the invalidation tags. The primary reason for these invalidation streams is to

²In this case in the form of an Object-Relational Mapper

³<http://www.ibm.com/analytics/us/en/technology/db2/>

allow transactional consistency such that the primary storage and cache nodes share information that can be used to fetch a consistent set of data across both sources.

The same approach also suggests using transactions between the primary storage and the application during the computation of cached objects to capture dependencies to entities from the primary storage using the queries made in the given transaction. This way the dependencies are captured automatically, but it requires the database to implement these transactions.

[Was11] suggests using static analysis of the code to capture dependencies between the cached functions and the underlying data. In this approach the programmer must denote the data relationships, which are analyzed by the static analyzer. The static analyzer will then detect relationships between the cached functions and the underlying data as well as between different cached functions. Based on this information an object in the application process will store the dependencies for each type of dependency and not for each dependency instance. The exact dependencies are then derived using queries to the primary storage when changes are triggered. The given approach also uses a static analyzer to detect the relevant triggers.

The advantages of the deploy-time approach is that it doesn't require additional processes and it does not have to keep track of state in the form of dependencies between the different cache object instances, but these decisions have the cost that it is difficult to implement fault-tolerant measures (such as invalidation retry) as well as a poorer user experience from the performance impact of invalidation. Furthermore the approach has no mechanisms for avoiding the problem of concurrency bugs (describe in section 3.6).

On the other end of the granularity scale, [CLT06] suggests a system that caches HTTP responses. It uses a sniffer process that monitors the lifetime of a HTTP request with the queries made to the storage system. Through the information captured by the sniffer, the system builds a table that maps a given HTTP resource to the queries made. The system then caches the HTTP resource that is invalidated when underlying data related to the given resource changes. This method is interesting since it allows to cache without changing the code of the web application, but it is only described at the granularity of HTTP responses since it uses the communication between the web application, storage system and cache to achieve automatic invalidation.

To be able to evaluate the techniques used in the automatic invalidation approaches described, we will divide the process into different sub problems. Automatic invalidation are based on a caching system that reacts and invalidates when changes happens to underlying data. Automatic invalidation are therefore

mainly working around requests, where the client requests to update underlying data in the primary store. If we consider figure 3.9 that illustrates this flow, the invalidation mechanism is responsible for step 3, 4 and 5. This involves the tasks of: *triggering cache invalidation* and *managing dependencies between underlying data and cached objects*. The following sections will consider these tasks and compare existing techniques.

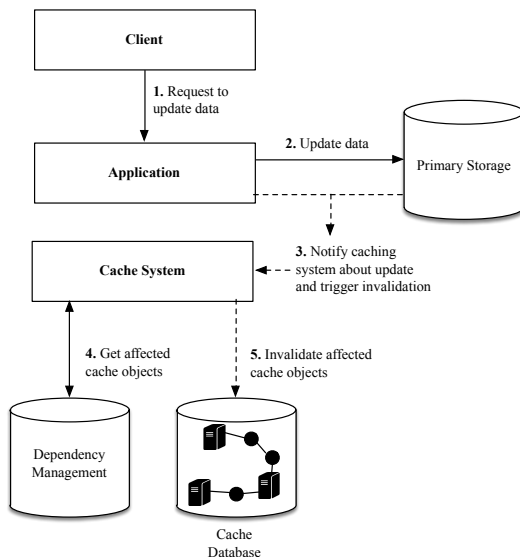


Figure 3.9: The control flow of automatic invalidation when a client requests to update underlying data

3.1.6.1 Triggering Cache Invalidation

When the client updates underlying data, the application receives a request from the client and sends a request to the primary storage to update the relevant data. To be able to react to this action, triggers have to be implemented during this flow. One technique used by [GZM11, Was11] is to have the triggers in the application such that the application triggers invalidation when the changes from the client has been applied to the primary storage. The given techniques involves a database wrapper (in the form of an ORM) that implements callbacks invoked when it sends commands to the database. If the database wrapper is implemented using the adapter-pattern, it is possible to change the underlying database technology without changing the implementation of the caching system. This means it is easier to implement the caching system for multiple

database technologies as well as change the technology for existing solutions.

Other solutions such as [Por12, CDI98, CID99] assumes that the database is able to send notifications when changes are made to the database. These information from the notifications are intercepted by the caching system and converted to invalidations. [Por12] uses an “invalidation stream” that is replicated directly across all cache nodes, which means the cache nodes has the responsibility of invalidating the correct cached objects. In [CDI98, CID99] this responsibility is extracted into a third process that converts the change notifications to invalidation of affected cache objects. A similar approach is used in [CLT06] that uses a proxy between the application and the primary storage to “sniff” the database traffic and relate it to the HTTP-request of the client.

In the techniques with in-application triggers, changes to the database around the application are not captured. So if the system has the requirement that the primary storage can receive commands from multiple applications, the best solution would be to use triggers directly from the database. But this also means the system is required to use a database technology that supports this. This comparison is also shown on figure 3.10.

	Advantages	Disadvantages
Database	- Any change is captured	- Require database to support triggers
Database Sniffer	- Any change is captured	- Require a database sniffer for the database technology used
Database Wrapper	- Supports all database technologies supported by the database wrapper or the API used by the database wrapper	- Changes made to the database around the database wrapper are not detected

Figure 3.10: Comparison of triggers for automatic invalidation

3.1.6.2 Dependency Management

After the invalidation has been triggered, the invalidation system needs to locate the cache objects that needs to be invalidated. This involves the task of identifying and declaring dependencies such that the triggers will invalidate the affected cache objects.

Since dependency management is a burden for the programmer and affects the correctness of the cache implementation, it would be most desirable to have fully automatic dependency management, which is achieved in [Por12, CLT06]. In [CLT06] use the technique of proxies between the different servers to sniff the traffic and thereby automatically derive dependencies between HTTP-responses and queries made during the request. [Por12] runs a transaction with the database while the cached object is compute and uses information from the

queries made during the transactions to derive dependencies between the underlying data and the cached object. These techniques removes the burden of cache management by having fully transparent caching, but it also means the programmer has less flexibility. Furthermore they are tightly coupled to the technologies used and makes it difficult to port the solutions to other technologies.

Other solutions relies such as [GZM11, Was11] on the programmer declaring dependencies from the cached objects to underlying data. Since the trigger can include information about which underlying data are changed, the caching system can use the declared dependencies to invalidate the corresponding cached objects. [GZM11] supports declarations through function calls that are stored in the memory of the application. The deploy-time model suggested in [Was11] uses static analysis of comments to allow the functions to be executed without the cache database. These techniques does not remove the burden of cache management completely, but they allow the programmer to specify the dependencies in a more declarative and robust way compared to using manual invalidation triggers.

The solution suggested by Jim Challenger et.al. [CDI98, CID99] uses dependencies declared in the content to construct an advanced dependency graph. When updates are made to content or underlying data, the affected cache objects are derived using the dependency graph. This solution is developed for a content management system, where the users can declare dependencies between the fragments of the content i.e. the dependencies are declared in each entity. This makes it unfeasible to use in application with slightly advanced data models, but it solves the problem well in the given case.

An overview of this discussion can be found in figure 3.11.

	Advantages	Disadvantages
Declared in code (Cache-Genie)	- Easy to reason about dependencies	- Relies on the developer correctly identifying and declaring dependencies
Declared in content (IBM)	- Allow different data source for different entities - The developer does not have to declare dependencies	- Burden for users to define dependencies on each entity
Static Analysis and Code Generation (Deploy-Time Model)	- Semi automatic; only require definition of database relations	- Static analysis cannot be implemented correct in dynamic languages - Difficult to detect relational dependencies - Requires knowledge about static analysis to implement
Database Transactions with Invalidation Tags (TxCache)	- Fully automatic	- Require the database to implement the transactions

Figure 3.11: Comparison of dependency management techniques for automatic invalidation

3.1.7 HTTP Caching

3.2 Choosing the Right Caching Technique

In general there is no best or correct solutions - it's a matter of choosing the solution best suited in the given context depending on the web application and specific use case. We will therefore compare the different approaches based on the parameters described in section 2.2. This comparison is illustrated on figure 3.12.

	Consistency	Strict Freshness	In-place update	Always Immediate Response	No Invalidation Management	Adaptability
Arbitrary Content						
Expiration-based Invalidation						High
Key-based Invalidation						High
Manual Trigger-based Invalidation						High
Key-based Invalidation with Async Update						High
Write-through Invalidation						Medium
Dan Ports: Transactional Consistency						Low
Chris Wasik: Managing Cache Cons...						Medium
Declared HTML-content						
IBM						Low
The other XML ones						Low
HTTP-response						
Y-K. Chang et al.: DB Sniffer						Low
DB-Queries						
Cache-Genie						Medium
Materialized Views						Medium

Figure 3.12: Comparison of caching approaches

From this comparison we see that some decisions have to be made between some of the parameters. At first the approach cannot both give an immediate response and provide strict freshness since if the cached computation takes Δt_c to compute and the invalidation happens just before the request then the response time will be Δt_c .

CHAPTER 4

Smache: Cachable Functions

The goal of this thesis is to present a caching solution that is able to handle long running computations by presenting content to the user fast while addressing the challenges of cache management and efficient update propagation. Based on the overview of existing caching solutions introduced in the last chapter, we will present a solution that solves the problem for the context of this thesis.

4.1 Existing Approaches are Not Sufficient

In the primary use case of the context in this thesis (described in section 1.3), the platform presents statistical information based on some advanced computation that takes a long time to compute (> 10 sec.). If we want to keep the teacher's attention to the platform, we need to find a caching technique that is not depending on the time taken to compute the information. Furthermore we want to keep the information as fresh as possible.

If we consider the overview on figure 3.12, we can already leave out caches that optimize DB-queries and declared content since they do not solve the problem of caching advanced statistical computations. The solution proposed by IBM

is made for declared HTML-content and not on computations based on data from a storage system. This leaves the set of approaches that are able to cache arbitrary content and HTTP-responses, but given caching arbitrary content has less assumptions, they are more attractive.

From the approaching caching arbitrary content it is desirable to have an approach with high adaptability, which was stated as one of the requirements of the system. From the approaches with high adaptability that has a response time that does not depend on the time of computation, are *expiration-based* and *manual trigger-based invalidation with asynchronous updates*.

From these approaches the expiration-based invalidation technique has the advantage that it doesn't require invalidation management, but it has the limitation that all cached values of the same kind are invalidated and updated at the same frequency.

If we consider use case example 4.1, the statistical information of current assignments are updated at the same frequency as the closed assignments, and if we want to keep the cached values for the current information up to date, we also need to update all non current information. This is not desirable with relation to efficiency since the CPU will be busy in time proportional to the time of the computing the statistical information for the assignment and the total number of assignments on the platform. In other words the number of CPU's we need to occupy for this job can be calculated by $\frac{\Delta t_c \cdot n_c \cdot u_c}{t_d}$, where Δt_c is the time of the computation, n_c is the number of computations, Δu_c is the number of updates per 24 hours and Δt_d is the number of seconds per 24 hours. Considering the example - if we have 500 assignments on the platform and we want to update the information every 10 minutes with a computation time of 30 seconds each, then we occupy a CPU in $\frac{60 \cdot 500 \cdot \frac{10 \cdot 60}{60 \cdot 60 \cdot 24}}{60 \cdot 60 \cdot 24}$ occupied CPU's = 50 occupied CPU's. If the amount of computational power isn't a problem, the expiration-based approach would be the best solution, but that is not the case of this thesis, where efficiency is a requirement.

The alternative with high adaptability is to use manual trigger-based invalidation with asynchronous updates after request as described in section 3.1.4. This approach has the advantage that the values are only invalidated (and thereby recomputed) when a trigger is invoked, which means the programmer has the opportunity to optimize the approach to only invalidate when it is relevant to the user e.g. when the cached value are supposed to change. Although this can be seen as an opportunity it can also be seen as a burden for the programmer to maintain these manual triggers. And by having asynchronous invalidations, the user is presented with a stale value, and only if the value is requested again after the update, the user receives the fresh value. If we consider example 4.1

it is not unlikely that the teacher looks at an old assignment some time after it has been closed. In cases such as this, where the value is only fetched once in a while it would not be optimal to have asynchronous updates, since the value will never be close to fresh when it is actually presented.

EXAMPLE 4.1 *On the Peergrade.io platform the teacher is presented with statistical information about the grades given by the students for a given assignment. These assignments are often current at specific time interval after which the assignment becomes “closed” and the statistical information are not updated.*

The approaches suggested by Chris Wasik and Dan Ports focus on automatic invalidation and could probably be extended to do asynchronous updates on request to allow immediate responses by giving up the option of freshness, which would give the same guarantees as the trigger-based invalidation with asynchronous updates with additional automatic invalidation, but it would also have the same problems.

The approach best fit for the context of this thesis is to use write-through invalidation that guarantees immediate response time as well as updating the content in-place such that the content cannot become older than the time taken to compute the value.

Although this is the best fit, it still leaves a burden for the programmer to declare the triggers that invokes the write through updates. Furthermore it leaves the challenge of ensuring the integrity of the cached values and invalidation marks in concurrent environments as described on figure 3.6.

The goal of this thesis is to present a caching approach that solves the problem of the thesis for the use case described in section 1.3 as well as fulfill the requirements described in section 1.2. Since the computation time in the given use case can be long, we must ensure that the response time does not depend on the computation time. Since we must respond immediately we need to trade-off strict freshness. To keep the values as up to date as possible we can use in-place updates based on trigger invalidation. As we can see on figure 3.12, there exists no solutions fulfilling those criteria. In the rest of the thesis we will present the design and implementation of the solution that meets these requirements, which we call “Smache”.

4.2 The Cachable Function Model

Smache uses the same programming model as the one described by Dan Ports [Por12]. The model is organized around *cachable functions*, which essentially are normal functions, where the result are *memoized*¹. The cachable functions are allowed to make request to the primary storage to fetch the underlying data for the function as well as calling other cached functions. This allows caching at different granularities that gives the benefit of optimizing the invalidation for different types of content.

4.2.1 Restricted to Pure Functions

As mentioned by Dan Ports [Por12] not all functions are cachable. To be able to cache the functions they need to be *pure*, which Dan Ports defines as: “[...] they [functions] must be deterministic, not have side effects, and depend only on their arguments and the storage layer state.”. Smache does not detect these properties in a function and therefore relies on the programmer to ensure only pure functions are cached.

4.2.2 Making Functions Cachable

The solution suggested by Dan Ports et.al. [Por12] includes fully automatic invalidation. The solution suggested by this thesis achieves a semi-automatic invalidation based on dependencies declared in the procedure that makes the function cachable. This thesis will extend from Dan Ports’ solution for making cached functions with arguments for declaring dependencies. The API is defined as following:

MAKE-CACHABLE(*fn*, *input-types*, *relation-deps*) → *cached-fn*: Makes a function cacheable. *cached-fn* is a new function that first checks the cache for the result of another call with the same arguments. If not found, it executes *fn* and stores its result in the cache.

The added arguments for the procedure does not alter the behaviour of the *cached-fn*, but it requires additional arguments used by the application for naming and invalidation. The *input-types* are the types of the arguments given to the original *fn* and the *relation-deps* are declared dependencies to data sources,

¹In our case: storing the result of the function and return the cached value when the function is called again with the same input

which cannot be derived from the input of *fn*. The *input-types* are used to indicate how the arguments should be interpreted with relation to naming (described in section 4.2.3 and used together with the *relation-deps* to achieve automatic invalidation (described in chapter 5).

4.2.3 Cache Object Localization

When the application fetches and stores the content of cached objects it need to have consistent naming such that the correct cached objects are fetched and the objects does not overwrite each other when stored. To make it easier to use cachable functions this is automatically handled by Smache. The name of a given cached object computed from a cached function *must be unique* such that no other cached function is able to have the same name.

The details on how this is achieved is specific to the implementation, but as an example, the Smache Python implementation uses a combination of the function name and the module in which it was defined. This becomes unique since you cannot use two modules of the same name.

4.2.4 Automatic Cache Invalidation

In Dan Ports' solution the system handles all parts of cache management such that the programmer only have to define which functions needs to be cached and not define how it should be invalidated. This is a great advantage since it avoids potential bugs related to cache invalidation, but it is a trade-off for flexibility as well as adaptability since this approach moves the complexity to the database in the form of assumptions about the primary storage, cache database as well as a daemon process for managing snapshots [Por12].

Smache does not remove the burden of cache management, but it will handle naming (or localization) and invalidation by only relying on the programmer declaring dependencies to underlying data. This will improve the usability of cache management since the cache invalidation only changes when there are changes to the underlying data as described more in chapter 5.

4.2.5 Data Update Propagation

Smache also offers the possibility to perform in-place updates using a data update propagation system inspired by the solution IBM used to achieve a 100% cache hit rate on the content management system for the Olympic Games in 2000 [CDI98, CID99].

Having in-place updates allows the programmer to have immediate response times while ensuring that the values returned are updated as soon as they are invalidated such that the value is not more stale than the time taken to compute it. This will be explained more in chapter 6.

4.3 Implementing Cachable Functions in Python

The basic cachable functions as they are described so far can be implemented using the architecture described in section 2.2. Smache implements the cachable function by exposing an implementation of the **MAKE-CACHABLE**. To apply the procedure the Smache library must be loaded in the application and applied. After the **MAKE-CACHABLE** function has been used, the control flow illustrated on figure 4.1 is applied. When a client makes a request involving a cached function, the application will call the function wrapper that tries to find the value in the cache first and if it is not present in the cache, the original function will be called after which the result is cached and returned to the caller.

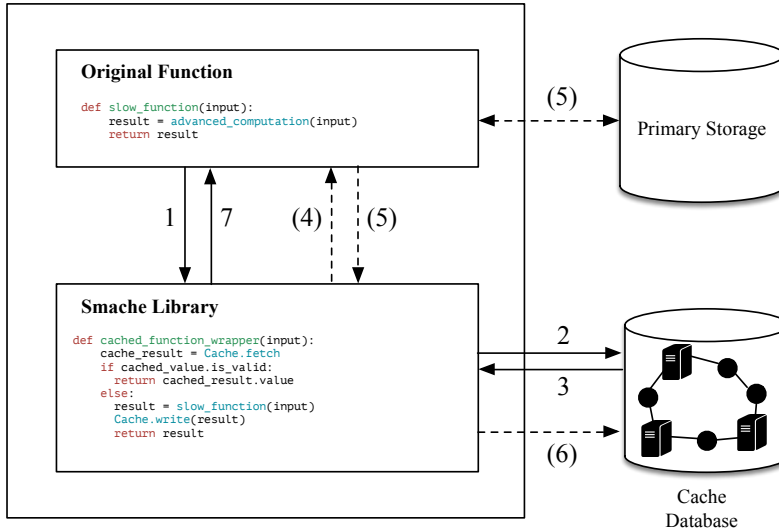
The procedure for making a function cached is implemented in using Python-decorators ² that is an annotation that can be used to modify the behaviour of existing functions while being able to keep a reference to the original function.

The implementation of the basic caching on the running example is seen on code snippet 4.2. In this code the only addition we have made is a call to a function called `computed` on a module `smache` with a single argument and a `@` as prefix that indicate it's a decorator. Additionally there is a decorator called `relation` in which we indicate the entity type as well as a function that describes how the given relation entity relates to an entity indicated through the input. This call corresponds to the implementation of **MAKE-CACHABLE**.

Smache requires the programmer to annotate the types of “data source”, which are passed into the original function. In this case we have annotated using the `Course`-class that corresponds to an object for an Object-Relational Mapper

²<https://www.python.org/dev/peps/pep-0318/>

Application



1. The application calls Smache through the cached function wrapper
2. Smache sends a **Lookup** request to the Cache Database
3. The cache adapter returns the value with an indication of freshness based on the result from the cache database
- (4). Smache calls the original function through the application
- (5). The application executes the original function, makes the necessary calls to the primary storage and returns the result
- (6). Smache sends a **Store** request to the Cache Database with the result from the application
7. Smache returns the value to the caller from the application

Steps annotated with **parenthesis and dashed lines** are only executed in case of **cache miss or cache failure**

Figure 4.1: The control flow during a call to a function cached through Smache

(ORM) that represents the collection *course*. The “data source types” supported are ORM-objects as well as the “Raw”-type that represents primitive values (such as strings, numbers etc.). The number of type annotations has to be the same

number of arguments given to the original function.

Smache uses the type annotations to do smart serialization and deserialization when the name of the cached value have to be inferred. The raw input are represented as they are in a serialized form, but the data source inputs are serialized using their. When the function is then deserialized the Smache library will use the id of the data sources to load a fresh version of the ORM-object from the database. Another advantage is that the Smache now know that this cached function has to be invalidated when the corresponding course entity is updated, which will be explained in chapter 5.

```

1  @smache.relations(
2      (Participant, lambda participant: participant.courses),
3      (Grade, lambda grade: grade.assignment_course)
4  )
5  @smache.computed(Course)
6  def course_score(course)
7      participants = Database.find_all_participants_in_course(course)
8      total_score = 0
9      for participant in participants:
10         total_score += time_consuming_participant_score(participant)
11     return total_score / len(participants)
12
13 @smache.relations(
14     (Grade, lambda grade: grade.graded_participant)
15 )
16 @smache.computed(Participant)
17 def time_consuming_participant_score(participant):
18     grades = Database.find_all_grades_for_participant(participant)
19     return numpy.advanced_statistical_method(participant)

```

Figure 4.2: Implementation of basic caching on the running example

4.4 Discussion

CHAPTER 5

Automatic Cache Invalidation

In the model for cachable functions described in the previous chapter 4 we've only described how the system can make a function cachable such that Smache automatically stores and locates the cached object returned by the function. While this removes the burden of naming and localizing cached objects, it does not solve the hard problem of invalidation.

To make cache invalidation easier with Smache, we will extend this solution with an invalidation mechanism that automatically invalidates the cached objects based on dependencies declared by the programmer.

This chapter describes how the suggested system uses the declared dependencies to track updates to underlying data and invalidate cached objects correctly. The solution for automatic invalidation is described in three parts. First we introduce the data structures used to achieve fast and efficient invalidation (section 5.1). Second, the *dependency registration* describes how the cached objects are registered into the ODG. Finally, the *invalidation propagation algorithm* (section 5.4) describes how Smache invalidates affected cached objects when changes are made to the underlying data. The chapter ends with a discussion on how the suggested solution meets the criteria and requirements of the thesis.

5.1 Simple Object Dependence Graph

The Object Dependence Graph (ODG) was first described in a series of papers by IBM [CDI98, CID99] and was build to support the content management system build for the Olympic Games in 2000. The final content served to the user is build from HTML-fragments and HTML-pages editable by the users as well as underlying data periodically changing. To be able to serve the final documents fast, the system pre-generates the HTML pages such that the system doesn't have to generate the pages on each request. To do this efficiently, the system maintains dependencies between the different kinds of objects and updates depending objects when underlying data changes. The cache object pre-generation can also be described in terms of caching, where the system uses an in-place caching approach.

Jim Challenger et.al. [CID99] presents a simple and a generalized version of the ODG. The generalized version is described for a content management system and includes a number of enhancements compared to the simple model that is unnecessary for the purpose of this thesis. We will therefore use the simple ODG that is described as following:

ODG can be represented using a directed graph, where a vertex either represents underlying data or an object that is a pure transformation of underlying data. An edge from a vertex representing underlying data u to a vertex representing an object o denoted (u, o) indicates that a change to u also affects o . [CID99] gives the following constraints for the simple ODG:

- Each vertex representing underlying data does not have an incoming edge
- Each vertex representing an object does not have an outgoing edge
- All Vertices in the graph correspond to underlying data (nodes with incoming edges) or objects (nodes with an outgoing edge)
- None of the edges have weights associated with them

Figure 5.1 illustrates an instance of a simple ODG with four vertices of underlying data (u_1 , u_2 , u_3 and u_4), two vertices representing objects (o_1 and o_2), and the dependencies between them. Where we can see that when underlying data u_2 changes, the system must update both the cached object of o_1 and o_2 and when u_4 changes, it only needs to update o_2 .

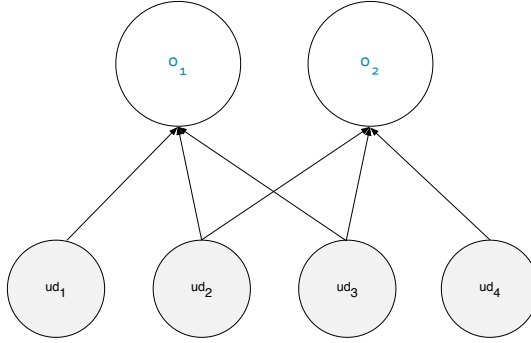


Figure 5.1: An example instance of a simple ODG

5.2 Dependency Data Structure for Cachable Functions

There exists two representations of the cachable functions. The first one is the static representation, which can be derived directly from the declarations in the source code. This representation only describes dependencies between the declarations and not cached object instances. We will define this representation as the *Declaration Dependence Graph (DDG)*. The other representation is dynamic and describes the dependencies between the entities of the underlying data and the cached object instances. We define this as the *Instance Dependence Graph (IDG)*.

5.2.1 Declaration Dependence Graph

When the cachable functions are defined as in section 4.2.2, the dependencies are declared using the entity types. The DDG is an extension of the Simple Object Dependence Graph, where the cached functions corresponds to the object vertices and the entity types corresponds to the underlying data. Where the simple ODG only includes direct dependencies, the DDG has two kind of dependencies. An edge from u to o denoted $(u, o)_d$ represents a direct dependency, which means that the dependency defines the given o . That is if a given cached object instance has a dependency to an instance of an entity then it cannot change. The direct dependencies includes an index indicating a position such that there is an order of direct dependencies for a given cached object. An edge from u to o denoted $(u, o)_l$ represents a lazy dependency, which means that the dependency can change throughout the lifetime of the cached object

instance and has to be evaluated before invalidation to derive the dependency for the instance. This could just as well be two graphs, but we model it as one for illustrative purposes.

Figure 5.2 shows the DDG for the running example from code snippet 1.1. For example we see that the cached object for the `participant_score`-function depends directly on the participant and it depends lazily on the grade. This is because there exists a participant for each score, but the dependencies from a cached object instance to any grade entity can change throughout the lifetime.

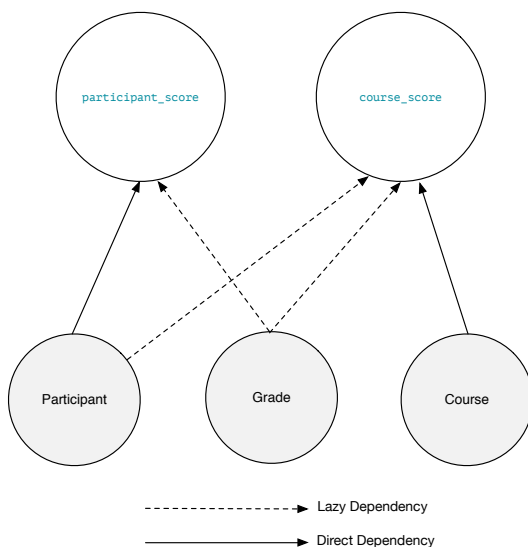


Figure 5.2: The Declaration Dependence Graph of the running example

The DDG graph is represented using two data structures. The lazy dependencies are represented using an outgoing adjacency list from the underlying data nodes. The direct dependencies are represented using an incoming adjacency list from the object nodes ordered by the position index of the dependency. To access a given adjacency list fast we use hash tables index by entity type for the outgoing adjacency list and by the object ID for the incoming adjacency lists as depicted on figure 5.3.

New dependencies can be added to the DDG using

- `add_dependency(s_id, t_id)`: Adds a dependency between the source id `s_id` used as index for the hash table and the target id `t_id` used as element in the adjacency list.

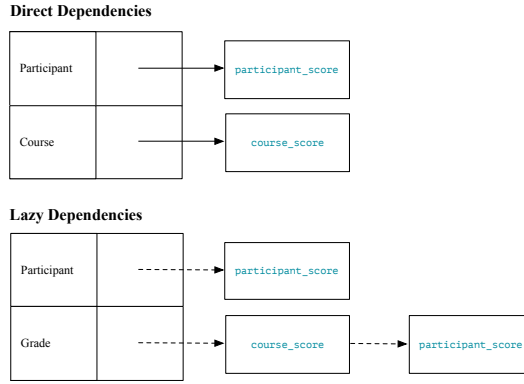


Figure 5.3: An illustration of the data structure representing the DDG on figure 5.2

When `add_dependency` is used we check if there already exists an adjacency list in the hash table using `s_id`. If there exists one, we add `t_id` to the given adjacency list, or else we add a new adjacency list only including `t_id` and indexed by `s_id`.

The DDG is build by iterating through all the cached functions with their respective dependencies and add the dependencies to the data structure as described above. The lazy dependencies are added with the entity type as `s_id` and the cached function as `t_id`. The direct dependencies are added using the cached function as `s_id` and the entity type as `s_id`. The advantage of the DDG is that it does not change in the lifetime of the application, which means it can be preprocessed when the application starts.

The queries needed for automatic cache invalidation are the following:

- `lookup_lazy_dependency(entity_type)`: Finds all lazy dependencies from a given entity type
- `lookup_direct_dependency(fun_id)`: Finds all direct dependencies from a given function id

Using these data structures we can access a pointer with access to all lazy or direct dependencies using $O(1)$ expected time if we use a hash table using perfect hashing. [FKS84]. In worst case the space used is the maximum number of edges $O(|u| \cdot |o|)$, where $|u|$ denotes the number of notes representing underlying data and $|o|$ represents the number of nodes representing objects.

5.2.2 Instance Dependence Graph

The Instance Dependence Graph (IDG) is also an extension of the simple Object Dependence Graph (ODG), where the data entities corresponds to underlying data and the cached object instance generated from the cachable functions corresponds to the objects. An edge from underlying data u to an object o denoted (u, o) indicates that if u is changed then o must be updated.

The instance dependence graph for the running example is illustrated on figure 5.4. In the given example, this graph is quite primitive. It's worth noting that the edges in the IDG are instance specific representations of the direct dependencies of the DDG seen on figure 5.2.

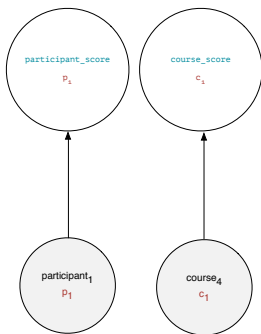


Figure 5.4: An example of an Instance Dependence Graph based on the running example

The representation of the IDG is similar to the representation of lazy dependencies for the DDG, where the dependencies are represented using an outgoing adjacency list indexed by a hash table as illustrated on figure 5.5.

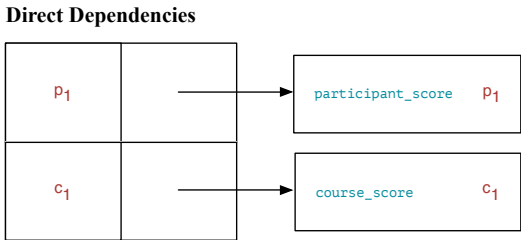


Figure 5.5: An illustration of the data structure representing the IDG on figure 5.4

New dependencies are added using `add_dependency` described for the DDG in section 5.2.1, where the entity id is the source and the cached object is the target.

The queries needed for automatic invalidation are:

- `lookup_cached_object(entity_id)`: Finds all cached objects depending on the `entity_id`.

The DDG uses the same amount of time and space as the DDG described in section 5.2.2.

5.3 Dependency Registration

For the application to know what and when to invalidate, the application must register the dependency declarations and name of cached functions instances. The dependency declarations that defines the lazy and direct dependencies are part of the source code and therefore available when the application is started as illustrated on figure 5.7. This registration flow simply collects all the cached functions registered through the cachable function procedure and adds them as dependencies in the DDG.

The registration of cached object instances happens when the cached object is accessed for the first time as illustrated on figure 5.6. In this flow we start by looking up the given cached object in the IDG, but since it's the first time the object is accessed it does not exist. To generate the name of the cached objects we use the ordered direct dependencies from the DDG combined with the input from the request. We derive dependencies to underlying data from the input that represents entities and add dependencies between the given entities and the cached object through the IDG.

Cached object requested for the first time

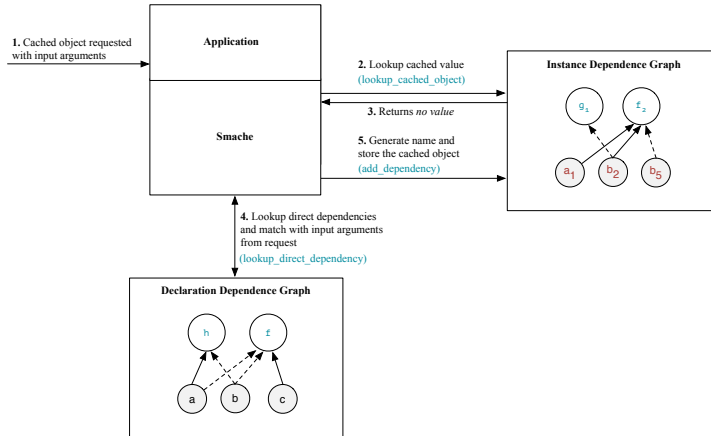


Figure 5.6: The flow in which cached object instances are accessed when they are accessed the first time

Application starts

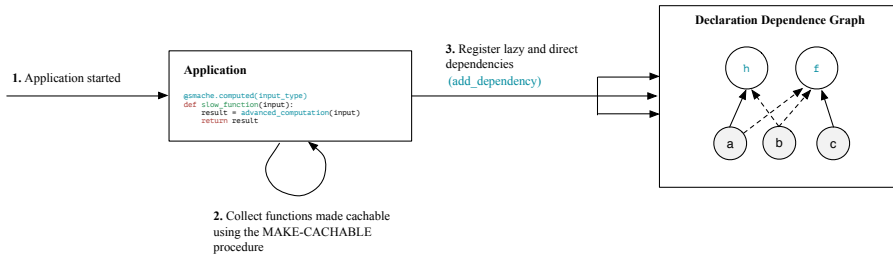


Figure 5.7: The flow in which lazy and direct dependencies are registered from the declarations

5.4 Invalidation Propagation

The purpose of invalidation is to be able to evaluate the freshness of a given cached object and know which cached objects are stale and need to be updated. A cached object is considered stale when it's underlying data has been updated, which happens during the request from a client that involves updating, inserting or deleting data in the primary storage. To be able to react to these events, Smache subscribes to callbacks from the database wrapper. When the database

transaction succeeds the database wrapper will notify Smache with the id and type¹ of the manipulated entity. This notification will trigger cache invalidation through Smache.

Based on the type and id of the changed entity, Smache updates affected cached objects in two steps:

1. Find keys for cached objects depending on the changed entity
2. Invalidate both set of keys using timestamp invalidation

There are two sets of object names found from the direct dependencies lazy and dependencies. To find the cached objects that directly depends on the given entity we query the IDG data structure with `lookup_cached_object(entity_id)`, which returns a list of names for cached objects affected by the given change. The other set of object names are found by evaluating lazy dependencies for the given entity type. The lazy dependency declarations are found by querying the DDG with `lookup_lazy_dependency(entity_type)`. The lazy dependency is then evaluated using the declared procedure, which returns a new set of entities directly related to the given data source. We then query the IDG with `lookup_cached_object(entity_id)` with id's from the new set of relations, which returns lists of all cached objects that depends on the given entities. From these cached objects we filter all the keys such that we only include keys for cached objects related to the lazy dependency.

We then invalidate both set of keys using timestamp invalidation as explained in the following section 5.4.1.

5.4.1 Timestamp Invalidation

To invalidate correctly the technique must comply with the liveness and safety property, which we define in the following terms:

- **Liveness:** *A stale object must eventually be invalidated*
- **Safety:** *The value representing a fresh object must be based on the newest version of the underlying data*

¹Type is another word for the relation in a relational database and a collection in a document-oriented database

In trigger-based invalidation where the name of a cached object is considered static we need external information to evaluate whether a given cached object is fresh. In a naive invalidation technique, we can use a boolean value that indicates whether or not the cached value is fresh. When invalidation is triggered the value is set to `false` and when the cached object has been updated it would be set to `true`. The problem here is that in an environment with multiple application servers, the technique would be prone to race conditions as illustrated on figure 3.6. Since the given cached object would be marked as fresh even though the value of the object in reality is stale, the cached object will not be invalidated until a new trigger is invoked and thereby contradict the liveness property.

One solution to this problem is to lock the invalidation mechanism such that the invalidation indication cannot be changed during an update. By using a lock the technique achieves liveness, but it also means that when another process wants to invalidate, it has to wait until the lock has been released. If the waiting process is a web server serving a user, the user would have to wait for the computation to finish.

To accommodate the liveness property and avoid user's having to wait for invalidation, we use a solution suggested in [CID99] that uses a form of logical timestamps to represent the version of a cached object. In this technique the cache system keeps track of the number of times a given cached object has been invalidated, which we denote *num_of_updates*. The number starts with 0 and is incremented every time a cached object's underlying data changes. We also keep a value representing the timestamp in which the current value of a cached object is based on, which we denote *last_update_timestamp*. When the cached object is recomputed, the *num_of_updates* is fetched and set as the timestamp of the current computation denoted *current_computation_timestamp*. When the computation returns the new value for the cached object, we fetch the value for *num_of_updates* again denoted *latest_num_of_updates*. At this point we have two cases:

1. If $latest_num_of_updates \geq current_computation_timestamp$ then we do nothing since a computation based on newer data has already updated the cached object.
2. Else we update *last_update_timestamp* to be equal to *current_computation_timestamp* and update the value.

TODO: Add a proof here that it has safety and liveness

5.4.2 Database Wrapper Triggers

We've already discussed existing approaches in section 3.1.6.1 that compares the techniques on figure 3.10. The triggers sent directly from the database or through a database sniffer have an advantage of capturing all changes made to the database - also those not sent through the application. The disadvantage of those techniques is that they are coupled to the database technology used. Instead of using a technique that depends on the exact database technology, Smache uses the database wrapper to trigger invalidations. The database wrapper makes it easier to change the implementation and thereby the database technology, which means the caching system becomes more flexible and easier to use in applications using different technologies. Furthermore since the triggers are intercepted through function callbacks in the web server process, the system doesn't need an external process to convert triggers from the database or sniffer into invalidations.

5.5 Implementing Automatic Invalidation

TODO: Write about how the automatic invalidation is implemented

5.6 Discussion

TODO: Discuss the following:

- How automatic invalidation meets the requirements
- Discuss how deriving dependencies instead of tracking would affect the solution (maybe do this in the same section as describing dependency tracking)
-

CHAPTER 6

Data Update Propagation

Until now we've explained how to build a caching system that is able to cache the result of functions, where the result is invalidated automatically when its underlying data changes. This solution makes it easier to manage cache invalidation and locate cached values, but after a result has been invalidated the user has to wait for the value to be recomputed, which can be critical for the user experience in cases where the computation time is too long. In this chapter we will extend the current solution with in-place updates using a data update propagation (DUP) algorithm that schedules updates for invalidated cache objects addressing the second challenge of the problem description 1.1. We will start by covering existing approaches (section 6.1) followed by an analysis of the problems involved with concurrent in-place updates (section 6.2). Based on the knowledge from these sections we will describe the design (section 6.3) and implementation (section 6.4) used in Smache.

6.1 Existing Data Update Propagation Approaches

In the approach suggested by Jim Challenger et.al. [CDI98, CID99, CIW⁺00], the DUP algorithm is based on invalidation from the Object Dependence Graph described in section 5.1. When underlying data changes all depending cached

objects are scheduled for update, but since some of the cached objects depends on other cached objects they cannot be computed in any order. If a cached object o_2 that depends on another cached object o_1 where updated first then it would be based on an old version of o_1 , which means it would still be stale and the algorithm would not have liveness. To solve this problem the approach updates all cached objects in a topological order, which ensures that o_2 is ordered after o_1 since it depends on o_1 . One limitations of this technique is that it only works if the object dependence graph has no cycles i.e. it is an *Directed Asyclic Graph*. Another limitation is that when there is an order of the jobs then they must be synchronized when executed in parallel.

TODOS:

- Write about Labrinidis's strategy
- Write about DBProxy

6.2 Race Condition on Write-Through Invalidation

When the name of the cached objects are the same as in trigger-based invalidation, the cached object becomes a shared resources that have to be accessed and updated from multiple web servers. This means we have to consider the concurrency challenges related to this concurrent environment to avoid avoid race conditions(while ensuring liveness of the system). As illustrated on figure 6.1 a race condition can occur when there are processes updating the same cached objects at the same time. In this case a race condition affects the correctness of the system such that a given cached object is marked as fresh even though it's value is stale.

6.3 The Data Update Propagation Algorithm

The existing solution described in section 6.1 assumes that a given cached computation is only accessed by one process at the same time. This has the advantage that they avoid race conditions, but it also means that the scheduling have to synchronize parallel execution to avoid updating the same object. The optimizations are then achieved using scheduling algorithms that evaluates the update jobs and prioritize them based on some metric such as the freshness or computation time. By having this assumption and enforce an order of execution

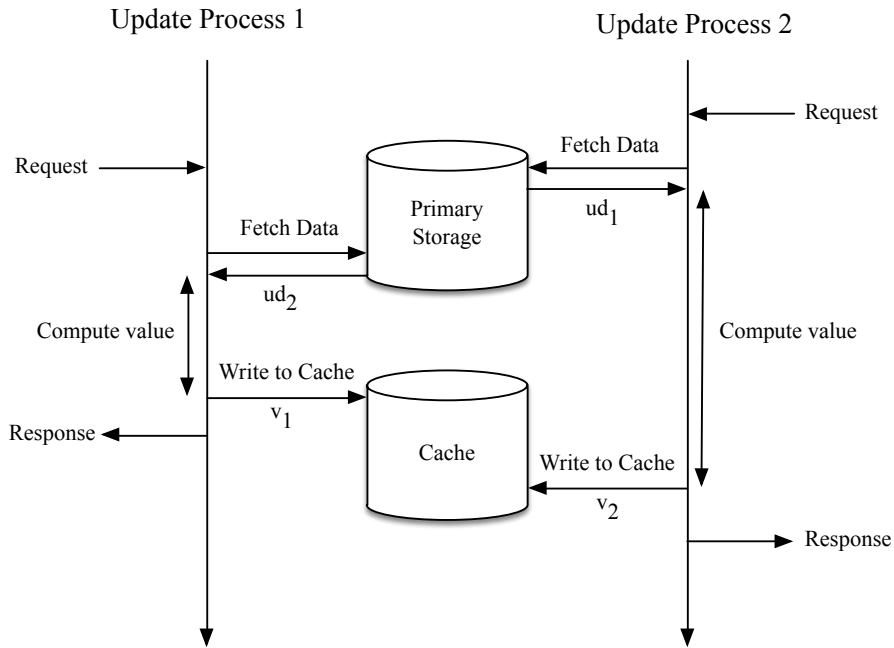


Figure 6.1: Showing how two concurrent caching updates from two different application servers results in an inconsistent state. We see that even though the request from *Update Process 2* are based on data older than *Update Process 1* it gets to write

some jobs have to wait to be executed. The only way to optimize the execution of these jobs is to acquire faster CPU's, which is expensive compared to buying more CPU's.

Instead of using advanced scheduling algorithms to achieve a high throughput we suggest a solution that allows concurrent updates of the same cached object. By allowing this we can scale our update execution by buying more CPU's and thereby be able to execute more jobs at the same time.

To allow concurrent updates we use Timestamp Invalidation as already explained in section 5.4.1. Timestamp Invalidation ensures that a given computation does not overwrite the value of a cached object if the computation is based on a newer version of underlying data. Figure 6.2

Since invalidation timestamps ensures the integrity of the cached objects we are

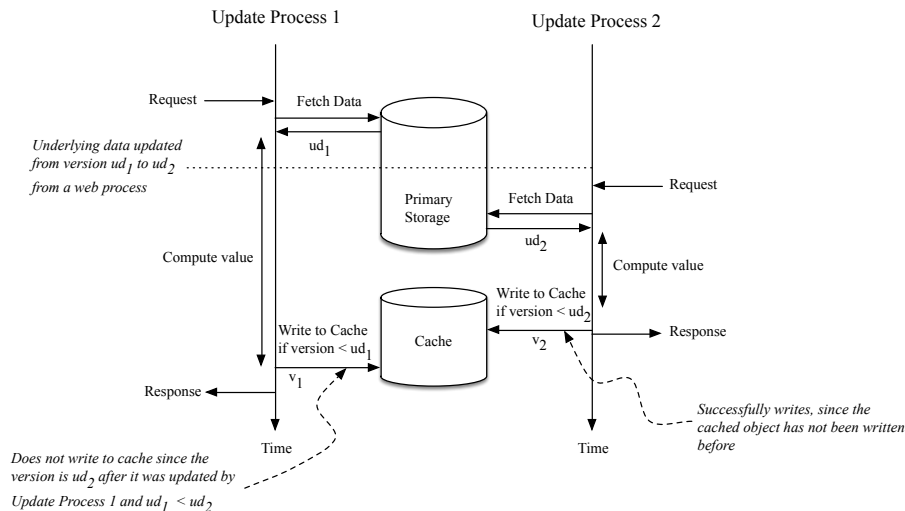


Figure 6.2: How Invalidation Timestamps fixes the concurrency problem described in figure 6.1.

allowed to scale horizontally and thereby execute as many jobs at once as there are update processes available. Furthermore since timestamp invalidation makes it possible to evaluate if a cached object is fresh, we can avoid computation of duplicates by only execute computations for stale cached objects. We could also optimize it even further by avoid scheduling recomputations for cached objects that are already scheduled, but this would require an atomic procedure to check if there already exists a similar job and only schedule the new job if there are none.

6.4 Implementing the Data Update Propagation Algorithm

TODO: Write about how the data update propagation is implemented

- Using queues (i.e. the producer/consumer model) as the scheduling mechanism (add job to queue to schedule update). Workers are then polling from the queue and executing.
- How the timestamp mechanism works (if it has not already been shown)

in the chapter on automatic invalidation)

- The architecture of the solution

CHAPTER 7

Results and Evaluation

Possible tests (choose the ones that makes sense):

- Prove that update requests are not made slower (because they are async)
- Test that functions are not tested when they are already fresh
- Test how often cached functions are executed unnecessarily
- Test how long it takes to warm up cache
- Test memory consumption of cache database when using Peergrade.io
- Test and compare sequential topologically sorted with parallelized
- Test update throughput under different write-patterns
- Test some QoD measurements

Discuss this:

- Software Design
- Adaptability
- (Fault-Tolerance)

CHAPTER 8

Conclusion

- Pull based more fault tolerant - if execution fails, we just serve an allright value

8.1 Future Work

APPENDIX A

Stuff

This appendix is full of stuff ...

A.1 Code Snippet for Trigger-based Invalidation with Asynchronous Update

```
1  def time_consuming_participant_score(participant):
2      return numpy.advanced_statistical_method(participant)
3
4  def cache_key_for_participant_score(participant):
5      cache_key_components = [
6          'cached_participant_score',
7          participant.type,
8          participant.id
9      ]
10     return '/'.join(cache_key_components)
11
12 def cached_time_consuming_function(participant):
13     cache_key = cache_key_for_participant_score(participant)
14
15     if is_fresh_in_cache(cache_key):
16         return fetch_from_cache(cache_key)
17     else:
18         update_cache_async(
19             time_consuming_participant_score,
20             participant
21         )
22     return fetch_from_cache(cache_key)
23
24 # Load the participant from the primary storage
25 participant = ParticipantDB.load_one_from_database
26
27 # Call the cached version of the time_consuming_participant_score
28 # Since there is currently no value in the cache it returns nothing
29 print cached_participant_score(participant)
30
31 # sleep for 5 seconds to wait for the computation to finish
32 sleep(5)
33
34 # This time we get an immediate response since the result is
35 # cached from the asynchronous update
36 print cached_participant_score(participant)
37
38 # Now we invalidate the cached value
39 cache_key = cache_key_for_participant_score(participant)
```

A.1 Code Snippet for Trigger-based Invalidation with Asynchronous Updates

```
40 invalidate_cached_value(cache_key)
41
42 # This time we serve the same value as when it was called
43 # previously, but now it is stale
44 print cached_participant_score(participant)
```


Bibliography

- [CDI98] J. Challenger, P. Dantzig, and A. Iyengar. A scalable and highly available system for serving dynamic data at frequently accessed web sites. In *Supercomputing, 1998.SC98. IEEE/ACM Conference on*, pages 47–47, Nov 1998.
- [CID99] J. Challenger, A. Iyengar, and P. Dantzig. A scalable system for consistently caching dynamic web data. In *INFOCOM '99. Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, volume 1, pages 294–303 vol.1, Mar 1999.
- [CIW⁺00] J. Challenger, A. Iyengar, K. Witting, C. Ferstat, and P. Reed. A publishing system for efficiently creating dynamic web content. In *INFOCOM 2000. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, volume 2, pages 844–853 vol.2, 2000.
- [CLT06] Yeim-Kuan Chang, Yu-Ren Lin, and Yi-Wei Ting. Caching personalized and database-related dynamic web pages. In *2006 International Workshop on Networking, Architecture, and Storages (IWNAS'06)*, pages 5 pp.–, 2006.
- [CRM91] Stuart K. Card, George G. Robertson, and Jock D. Mackinlay. The information visualizer, an information workspace. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '91*, pages 181–186, New York, NY, USA, 1991. ACM.

- [FKS84] Michael L. Fredman, János Komlós, and Endre Szemerédi. Storing a sparse table with $O(1)$ worst case access time. *J. ACM*, 31(3):538–544, June 1984.
- [GZM11] Priya Gupta, Nickolai Zeldovich, and Samuel Madden. *Middleware 2011: ACM/IFIP/USENIX 12th International Middleware Conference, Lisbon, Portugal, December 12-16, 2011. Proceedings*, chapter A Trigger-Based Middleware Cache for ORMs, pages 329–349. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011.
- [Han12] David Heinemeier Hansson. How key-based cache expiration works, February 2012.
- [Joh10] Robert Johnson. More details on today’s outage. website, September 2010.
- [Mil68] Robert B. Miller. Response time in man-computer conversational transactions. In *Proceedings of the December 9-11, 1968, Fall Joint Computer Conference, Part I*, AFIPS ’68 (Fall, part I), pages 267–277, New York, NY, USA, 1968. ACM.
- [Mula] Open Source Multiple. Memcached protocol documentation.
- [Mulb] Open Source Multiple. Transactions in redis.
- [Por12] Dan R. K. Ports. *Application-Level Caching with Transactional Consistency*. Ph.D., MIT, Cambridge, MA, USA, June 2012.
- [Was11] Chris Wasik. Managing cache consistency to scale dynamic web systems. Master’s thesis, University of Waterloo, 2011.