

# Spis treści

<b>1</b>	<b>Wstęp</b>	<b>1</b>
1.1	Problematyka i zakres pracy . . . . .	2
1.2	Założenia wstępne . . . . .	3
1.3	Układ pracy . . . . .	4
<b>2</b>	<b>Analiza jakościowa i wydajnościowa aplikacji</b>	<b>5</b>
2.1	Analiza wydajności serwera WWW . . . . .	6
2.2	Analiza wydajności frontendu . . . . .	12
2.3	Analiza wydajności bazy danych . . . . .	15
2.3.1	Optymalizacja table przy użyciu indeksów . . . . .	18
<b>3</b>	<b>Wymagania i budowa aplikacji</b>	<b>24</b>
3.1	Zakres funkcjonalny . . . . .	24
3.2	Budowa aplikacji . . . . .	27
3.3	Projekt bazy danych . . . . .	29
3.4	Organizacja kodu . . . . .	29
<b>4</b>	<b>Architektura aplikacji</b>	<b>31</b>
4.1	Rodzaje chmur . . . . .	31
4.2	Nie relacyjne bazy danych NoSQL . . . . .	32
4.3	Aplikacja zorientowana na usługi . . . . .	32
<b>5</b>	<b>Optymalizacja aplikacji</b>	<b>33</b>
<b>6</b>	<b>Optymalizacja kodu klienta</b>	<b>34</b>
<b>7</b>	<b>Metody rozproszenia aplikacji i usług</b>	<b>35</b>
	<b>Bibliografia</b>	<b>36</b>
	<b>Spis rysunków</b>	<b>37</b>
	<b>Spis listingów</b>	<b>38</b>
	<b>Płyta CD</b>	<b>39</b>

## **Abstract**

Nowadays Internet has been used as a wide range marketing and business tool. Most people is using Internet in various fields of life. For this reason sites and web applications becomes crowded and overloaded. Some times they can even stop working which will result in loss of money. That is why creating of high performance and stable applications is so important today. Also giving possibility to access page by huge amount of users can be cost-effective, when we put some commercials.

Last but not least web application architecture needs to be considered. It is very common that companies are using expensive hosting solutions, that doesn't actually fit their needs. On the another hand, when periodically massive traffic comes, their servers cannot handle this. Thats why cloud service providers are considered to be a good alternative in this situations, because companies are paying only for resources actually used and nothing more.

This thesis will cover the topic of efficiency optimization of web applications. It will be splitted into 3 main areas of optimizations. Firstly optimization of application architecture will be considered and then other issues like frontend optimization as well as database and webserver tuning will be discussed. Also this thesis needs to provide some useful tools and techniques that need to be used for mesuring of actual application performance and availability. Without this knowledge we cannot even start the main part of optimization, because we do not find the typical bottlenecks.

# Rozdział 1

## Wstęp

Niniejsza praca dyplomowa dotyczy zagadnień inżynierii oprogramowania oraz technologii baz danych wykorzystanych w dziedzinie *e-commerce*. Główny cel badań stanowi przedstawienie realnych korzyści biznesowych wynikających z zastosowania szerokiego spektrum usprawnień aplikacji internetowych.

W treści pracy przedstawiona jest analiza technik optymalizacji aplikacji internetowych z uwzględnieniem wykonywanych po stronie serwera (*server side*) oraz szeregu usprawnień po stronie przeglądarki *client side*. Równie istotny jest wybór metod dających najlepsze rezultaty w świetle obecnych technologii. Dodatkowo ważne jest wyróżnienie wszystkich pośrednich czynników, które mogą w jakimkolwiek stopniu wpłynąć na działanie aplikacji.

Do implementacji systemu wykorzystano technologie skryptowe PHP oraz Python. Pierwsza z nich pozwoli na szybkie przedstawienie obrazu typowej aplikacji e-commerce (olbrzymi odsetek takich aplikacji w internecie jest napisanych właśnie w tym języku). Python z kolei pozwoli wykorzystać zalety platformy Google Application Engine (GAE), która zapewnia skalowalną architekturę do późniejszych testów.

Obserwacja zostanie przeprowadzona na przykładzie aplikacji z dziedziny *e-commerce* - księgarni elektronicznej. Wybór takiej tematyki jest celowy, ponieważ najczęściej właśnie w takich aplikacjach występują problemy natury optymalizacyjnej. Spowodowane jest to najczęściej koniecznością obsłużenia wielu klientów, transakcji bazodanowych, czy przede wszystkim generowanie rozbudowanego wizualnie interfejsu użytkownika. Podczas generowania obrazów, wykonywania kodu serwera, czy interpretowania rozbudowanych struktur dokumentu *HTML*, czas oczekiwania na wynik może się zauważalnie wydłużyć.

Praca ma również na celu prezentację narzędzi badawczych umożliwiających znalezienie wąskich gardeł aplikacji. Tylko sukcesywne łączenie różnych

narzędzi oraz ciągle monitorowanie działania zapewni oprogramowaniu stabilność oraz wysoką dostępność.

### 1.1 Problematyka i zakres pracy

Wraz ze wzrostem popularności Internetu jako medium informacyjnego, istotnym problemem stało się obsłużenie napływającego ruchu sieciowego ze strony użytkowników. Pojęcie czas to pieniądz ma tutaj kluczowe znaczenie, ponieważ umiejętność przetworzenia jak największej ilości użytkowników w jak najkrótszym czasie będzie przekładała się na realne zyski.

Innym ważnym problemem dotyczącym stron jest zapewnienie wysokiej dostępności usługi, czyli wyeliminowanie do minimum wszelkiego rodzaju przerw wynikających z błędów aplikacji. Temat ten związany jest jednak głównie z odpowiednią konfiguracją sprzętową czyli wykorzystaniem równoległe działających instancji sprzętowych, które będą wykorzystywane w celu zrównoważenia ruchu, lub awaryjnie w wypadku uszkodzenia nośnika danych na jednym z urządzeń.

W internecie istnieje stosunkowo dużo publikacji związanych z tematyką optymalizacji aplikacji webowych, jednakże w większości wypadków omawiany jest tylko nikły procent wszystkich zagadnień. Zazwyczaj pomijane są aspekty związane z kodem po stronie klienta, a także studium narzędzi i metod badawczych. Najpopularniejsze na rynku są publikacje dotyczące optymalizacji samego kodu lub zapytań bazodanowych w zależności od użytego języka aplikacji lub bazy danych.

Praca ma na celu przedstawienie możliwie najszerszego wachlarza technik optymalizacji. Należy przy tym zaznaczyć, że statystycznie tylko 20% czasu przetwarzania strony przez przeglądarkę, jest poświęcane na oczekiwanie na odpowiedź serwera. Implikuje to olbrzymie znaczenie optymalizacji kodu po stronie klienta w celu znacznego przyspieszenia odpowiedzi. Nie bez znaczenia są też czynniki takie jak lokalizacja geograficzna strony oraz konfiguracja serwera.

Obecnie Internet przestał już być tylko wojskowym eksperymentem, czy zaledwie miejscem na prezentację własnej strony domowej. Wyewoluował on do medium globalnej komunikacji z gigantyczną ilością klientów docelowych. Większość firm, instytucji czy organizacji rządowych czuje się w obowiązku posiadania i utrzymywania strony internetowej, zazwyczaj spełniającej określone cele biznesowe. Pokazuje to jak bardzo powszechnym narzędziem codziennego użytku jest dzisiaj globalna pajęczyna.

Ze względu na niekwestionowaną popularność języka PHP oraz jego ol-

brzymią prostotę - w stosunkowo krótkim czasie od powstania języka, zaczęły pojawiać się proste strony, następnie aplikacje internetowe a kończąc na portalach i usługach sieciowych. Język PHP stał się narzędziem na tyle uniwersalnym, że zagadnienia modelowane przy jego użyciu, można z powodzeniem przenieść na inne platformy takie jak *Java Enterprise* czy *.NET*. W sieci istnieje ponadto wiele gotowych implementacji systemów e-commerce: sklepów (np. Magento), systemów CRM (SugarCRM) czy np. platforma edukacyjna Moodle będąca częścią infrastruktury edukacyjnej Politechniki Łódzkiej np. dla Wydziału FTIMS. Wymienione przykłady zostały w całości zaimplementowane w języku PHP, a o ich popularności świadczą miliony ściągnięć i wdrażeń.

Jedną z wad gotowych rozwiązań jest fakt, że nie zawsze są one dostosowane do wszystkich stawianych przed projektem wymagań. Powoduje to, że, projekt docelowy w rezultacie otrzyma więcej funkcjonalności niż jest to wymagane. Z drugiej jednak strony możliwe jest, że gotowe rozwiązanie będzie wymagało szeregu usprawnień lub dodania nowych funkcjonalności. W większości wypadków, rozwiązania gotowe nie są jednak od początku dostosowane do bardziej zaawansowanych zastosowań lub wymagań wydajnościowych. Dlatego ważne jest by wykorzystując istniejące narzędzia skalować aplikacje do konkretnych potrzeb lub przewidzieć przyszłe obciążenie.

Projektowana w ramach pracy aplikacja stanowi studium przypadku analizy wydajnościowej aplikacji działającej w ściśle określonej architekturze. W ramach analizy omówione zostaną następujące zagadnienia:

- metody badawcze
- dobór odpowiedniej technologii,
- wybór właściwej architektury sprzętowej,
- projekt i optymalizacja bazy danych,
- optymalizacja kodu klienta,
- rozproszenie usług

### 1.2 Założenia wstępne

Treść pracy dyplomowej stanowi wypadkową informacji zawartych w dokumentacjach dotyczących użytych technologii, jak również wiedzy autora zdobytej podczas implementacji wielu zróżnicowanych aplikacji internetowych. Bardzo istotne dla publikacji były również informacje pochodzące ze sprawdzonych źródeł takich jak np. oficjalny blog programistyczny Yahoo. Serwis Yahoo ze względu na olbrzymie doświadczenie w kwestii skalowania aplikacji internetowych postanowił podzielić się tą wiedzą na łamach swoich

stron internetowych oraz kilku specjalistycznych książek.

W kolejnych rozdziałach omawiane będą kolejne etapy drogi, jaką pokonuje żądanie od rozpoczęcia do zwrócenia i zinterpretowania przez przeglądarkę użytkownika. Często kolejne etapy są prawie niezauważalne ze względu na małe różnice czasowe, jednak podczas wnikliwej analizy każdy z etapów trzeba rozpatrywać indywidualnie.

Analiza wydajnościowa przeprowadzana w wypadku prostych aplikacji, które nie będą w przyszłości podlegały intensywnemu obciążeniu nie ma w zasadzie sensu. Publikacja zyskuje na wartości w wypadku kiedy projekt jest bardziej rozbudowany, a my potrzebujemy szybkiego i rzetelnego sposobu na wykrycie potencjalnych elementów do optymalizacji.

### 1.3 Układ pracy

Praca została podzielona na następujące rozdziały:

- **Rozdział pierwszy** opisuje narzędzia badawcze, a także wyjaśnia teorię działania aplikacji opartych o protokół HTTP
- **W drugim rozdziale** przedstawiono wymagania stawiane przed aplikacją oraz jej budowę.
- **Trzeci rozdział** dotyczy optymalizacji związanych z wyborem odpowiedniej architektury.
- **Rozdział czwarty** dotyczy optymalizacji związanej z implementacją aplikacji
- **Rozdział piąty** dotyczy optymalizacji kodu klienckiego
- **Rozdział szósty** metody rozproszenia aplikacji i usług
- **Rozdział siódmy** podsumowanie i wnioski końcowe

## Rozdział 2

# Analiza jakościowa i wydajnościowa aplikacji

Projektując aplikacje, już w fazie projektowej, należy myśleć o zapewnieniu wysokiej wydajności, a także o potencjalnych problemach, które mogą się pojawić po wdrożeniu oprogramowania. W celu zapewnienia tworzonej aplikacji najwyższej skuteczności pracy, należy wziąć pod uwagę wiele cech, wśród których najważniejsze zdefiniowane są poniżej.

### **Skalowalność (ang. *Scalability*)**

cecha aplikacji, określana jako zdolność do wzrostu wydajności aplikacji wraz ze zwiększeniem ilości dostępnych zasobów sprzętowych (serwery WWW, bazy danych, wydajniejsze procesory).

### **Niezawodność (ang. *High availability*)**

stanowi projekt, jak i odpowiednią implementację systemu, zapewniającą określony poziom ciągłości wykonywania operacji w czasie. Polega to na zapewnieniu jak największej dostępności usługi.

### **Wydajność (ang. *Performance*)**

przekłada się na możliwość szybkiego wykonywania kodu aplikacji oraz utrzymaniu czasu odpowiedzi aplikacji na stosownym poziomie.

Często skalowalność jest mylona z wydajnością, jednak przekładając to na bardziej życiowy przykład, wydajność aplikacji można porównać do szybkiego samochodu. Z drugiej strony, bez zapewnienia odpowiednich dróg, ten szybki samochód lub ich grupa, nie jest w stanie rozwinąć maksymalnej prędkości. W najgorszym wypadku może nawet utknąć w korku, blokowany przez inne pojazdy. Skalowalność jest więc zapewnieniem **odpowiedniej infrastruktury** gwarantującej właściwy rozrost systemu.

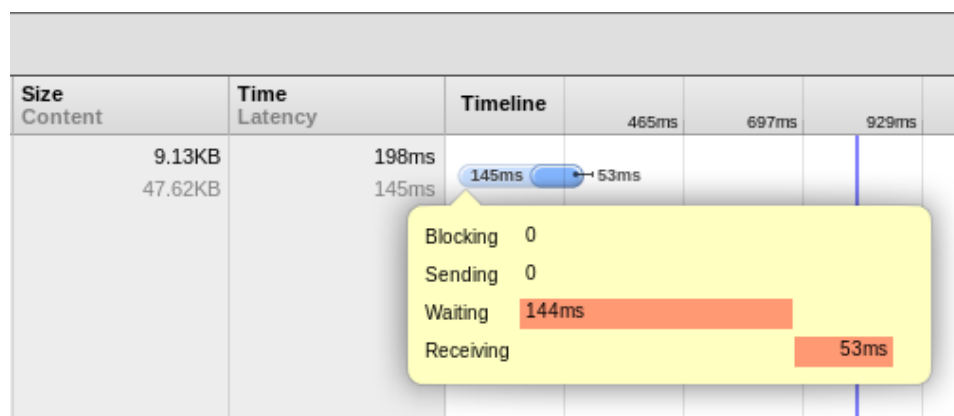
W celu zapewnienia możliwie najlepszej jakości tworzonej aplikacji, należy stale monitorować aktualny poziom wydajności aplikacji. Należy jednak

mieć na uwadze, że na wydajność aplikacji składa się czas przetwarzania poszczególnych węzłów systemu. Należy więc testować każde z nich osobnymi metodami, omówionymi w dalszej części pracy.

Szukając przyczyn błędów, warto rozpocząć analizę od najbardziej ogólnego komponentu czyli serwera WWW odpowiedzialnego za wysyłanie odpowiedzi na żądanie użytkownika. Następnie należy dokonać dekompozycji, wyróżniając kolejne węzły systemu, takie jak dalsze instancje serwera WWW czy serwery bazodanowe.

## 2.1 Analiza wydajności serwera WWW

Zadaniem serwera WWW jest wysłanie do inicjatora żądania wyniku przetwarzania zasobu określonego adresem URL. W najprostszym wypadku, analiza wydajności serwera, polega na odpytaniu go o określony zasób i zmierzenie czasu od rozpoczęcia tej akcji, do odebrania rezultatu. Taki proces można prześledzić i przeanalizować w większości popularnych przeglądarek np. Google Chrome, które jest wyposażone w wiele przydatnych narzędzi analitycznych (Rys. 2.1).



Rys. 2.1: Analiza czasu wykonywania strony <http://ftims.edu.p.lodz.pl/> wykonana w przeglądarce Google Chrome

Taki sposób analizy, jest jednak przydatny jedynie w wypadku znacznych problemów z wydajnością aplikacji, ponieważ testowanie czasu odpowiedzi dla pojedynczego użytkownika, wykonującego pojedyncze żądanie, nie jest w żadnym stopniu miarodajne.

W celu zapewnienia bardziej rzetelnego testu, należy skorzystać z dedykowanych rozwiązań takich jak **ab** oraz **siege**. Są to typowe narzędzia przeznaczone do sprawdzania jak dobrze serwer radzi sobie z obsługą bardziej złożonego ruchu sieciowego. Przykładowo, dla wcześniej użytej strony, można zasymulować ruch równy wykonaniu 10 jednoczesnych żądań przez



10 niezależnych użytkowników. W tym celu należy wydać komendę `ab -n 10 -c 10 http://ftims.edu.p.lodz.pl/`. Rezultat działania komendy widoczny jest na listingu 2.1.

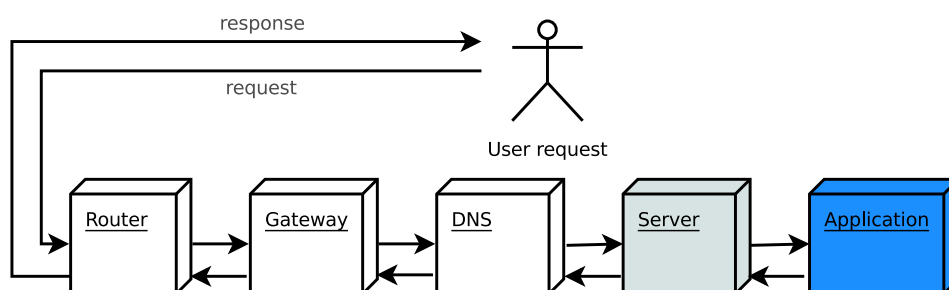
Listing 2.1: Analiza strony z wykorzystaniem narzędzia `ab`

```
1 Server Software:      Apache/2.2.14
2 Server Hostname:      ftims.edu.p.lodz.pl
3 Server Port:          80
4
5 Document Path:        /
6 Document Length:      48759 bytes
7
8 Concurrency Level:     10
9 Time taken for tests:  0.695 seconds
10 Complete requests:    10
11 Failed requests:      0
12 Write errors:         0
13 Total transferred:    492510 bytes
14 HTML transferred:    487590 bytes
15 Requests per second:  14.38 [#/sec] (mean)
16 Time per request:     695.270 [ms] (mean)
17 Time per request:     69.527 [ms] (mean, across all concurrent
    requests)
18 Transfer rate:        691.77 [Kbytes/sec] received
19
20 Connection Times (ms)
21      min    mean [+/-sd] median    max
22 Connect:    52     63   7.4      63      74
23 Processing: 300    503  90.9     532     632
24 Waiting:    127    233  62.8     235     381
25 Total:      354    565  92.9     593     695
26
27 Percentage of the requests served within a certain time (ms)
28  50%      593
29  66%      612
30  75%      625
31  80%      628
32  90%      695
33  95%      695
34  98%      695
35  99%      695
36 100%      695 (longest request)
```

Jak można wywnioskować z powyższych danych, narzędzie wykonuje wiele przydatnych analiz, a także wyświetla informacje o badanym zasobie. Widać przede wszystkim, że czas oczekiwania na stronę przy 10 użytkownikach jest prawie trzykrotnie dłuższy, niż podczas jednego żądania wykonanego w przeglądarce (2.1).

Oczywiście na wyniki pomiarów ma także wpływ prędkość połączenia internetowego, dlatego w celu pominięcia dodatkowych czynników, testy docelowej aplikacji będą wykonywane przede wszystkim na lokalnym serwerze. Najbardziej miarodajną jednostką określającą wydajność aplikacji w wypadku narzędzia `ab` jest liczba zapytań na sekundę (*req/s*). Określa ona maksymalną ilość żądań w jednostce czasu, jaką aplikacja jest w stanie obsłużyć. Oczywiście im większa wartość, tym lepsza ogólna wydolność aplikacji.

Na rysunku 2.2 przedstawiono cykl życia żądania od użytkownika je inicjującego, kończąc na odebraniu odpowiedzi serwera. Jak można zauważyć żądanie przebywa stosunkowo długą drogę, nim trafi do faktycznej aplikacji. Wynikiem tego są dodatkowe opóźnienia zależne od stopnia skomplikowania architektury trzech pierwszych węzłów. Dlatego też należy mieć na uwadze, że problemy z szybkością działania aplikacji nie muszą leżeć wyłącznie po stronie aplikacji lub serwera. Do najbardziej popularnych należą: niska przepustowość łącza internetowego klienta, wolny serwer DNS, daleka lokalizacja geograficzna serwera WWW, źle skonfigurowany router lub bardzo obciążona sieć lokalna.



Rys. 2.2: Cykl życia żądania

Nawiązując do listingu 2.1, wartościami związanymi ze wspomnianymi w poprzednim akapicie węzłami są **Connect** oraz **Waiting**, czyli odpowiednio czas oczekiwania na połączenie z zasobem i czas pobierania odpowiedzi z zasobu. Istnieje 5 głównych czynników wpływających na czas odpowiedzi serwera.

### Położenie geograficzne i problemy z siecią komputerową.

Nie bez znaczenia dla czasu odpowiedzi, jaki użytkownik odczuwa jest też lokalizacja serwerów stron. Jeśli serwery są zlokalizowane w USA, a użytkownicy odwiedzający stronę są np. z Europy, dystans jaki musi pokonać żądanie od momentu dotarcia do zasobu, oczekiwania, aż do jego pobrania jest nie współmiernie większy niż w wypadku stron hostowanych dla tego samego położenia geograficznego. Stopień opóźnienia jest zazwyczaj uzależniony od ilości routerów, serwerów pośrednich, a nawet oceanów, które pokonuje żądanie od punktu początkowego do odbiorcy i z powrotem.

### Wielkość dokumentu odpowiedzi serwera.

Zależność między wielkością dokumentu, a czasem odpowiedzi serwera jest oczywista, łatwo więc sprawdzić, że im większy dokument trzeba pobrać, tym więcej czasu potrzeba na zakończenie tego procesu.

### Wykonywanie kodu aplikacji.

Najczęstsza przyczyna wolnego działania aplikacji wynika właśnie z braku

optymalizacji kodu klienta. Długi czas wykonywania kodu aplikacji implikuje, długi czas łączny oczekiwania na odpowiedź serwera. Problem ten zostanie szczegółowiej poruszony w rozdziale 5.

### Rodzaj użytej przeglądarki.

Nie bez znaczenia dla ogólnego czasu ładowania strony jest również rodzaj użytej przeglądarki. Często wbudowane w przeglądarkę wewnętrzne mechanizmy buforowania zasobów pozwalają w znaczny sposób zredukować ilość zapytań wysyłanych do serwera. Dotyczy to zwłaszcza danych statycznych takich jak arkusze CSS, pliki JavaScript czy zasoby graficzne, które nie zmieniają się zbyt często.

### Konfiguracja serwera WWW.

W zależności od użytej technologii, istnieje wiele różnych serwerów HTTP. Wśród najczęściej używanych, prym wiodzie serwer HTTP *Apache*. Dla rozwiązań napisanych w technologii Java często wykorzystywane są serwery *GlassFish*, *Tomcat*, *Jetty*. W większości wypadków zaraz po instalacji, oprogramowanie serwera nie nadaje się jeszcze do wykorzystania w produkcji. Należy wyważyć ustawienia serwera do bieżących potrzeb, ponieważ w większości wypadków domyślne ustawienia mogą znacznie obniżyć ogólną wydajność. Innym ważnym działaniem jest dostosowywanie serwera do konkretnych zastosowań - do serwowania plików statycznych lepszym rozwiązaniem jest wykorzystanie bardziej oszczędnego pamięciowo i operacyjnie serwera *Ngnix*, natomiast do bardziej zaawansowanych zastosowań, w tym wykonywanie kodu aplikacji, serwera *Apache* lub osobnej instancji serwera *Ngnix*.

Administratorzy serwerów WWW mają bezpośredni dostęp do statystyk odwiedzalności stron, przez co pozwala to zaobserwować pewne trendy odwiedzin użytkowników. Często jest tak, że dane zasoby są dużo intensywniej odczytywane przez użytkowników np. w czasie 10 minut stronę odwiedza 100 użytkowników. Łatwo to sobie wyobrazić np. w wypadku premiery jakiejś nowej gry, lub publikacji wyników egzaminu na uczelni. Taki periodyczny, lecz bardzo wzmożony ruch może powodować pewne trudne do ustalenia problemy z działaniem aplikacji. Dlatego też twórcy narzędzia *ab*, zaimplementowali również możliwość testów czasowych (ang. *timed tests*). W ten sposób można zasymulować jak strona będzie się zachowywała również w takich nagłych wypadkach.

Wydając komendę `ab -c 10 -t 30 http://ftims.edu.p.lodz.pl/`, można sprawdzić, jak zachowa się aplikacja odwiedzana przez 10 użytkowników jednocześnie w czasie 30 sekund. Ta komenda pozbawiona jest parametru *-t ilość żądań*, oznacza to że symulacja zakończy się po 30 sekundach lub po osiągnięciu limitu 50 000 żądań.

Listing 2.2: Test obciążenia czasowego

```
1 Benchmarking ftims.edu.p.lodz.pl (be patient)
2 Finished 504 requests
3
4 Server Software:      Apache/2.2.14
5 Server Hostname:      ftims.edu.p.lodz.pl
6 Server Port:          80
7
8 Document Path:        /
9 Document Length:      48759 bytes
10
11 Concurrency Level:    10
12 Time taken for tests:  40.180 seconds
13 Complete requests:    504
14 Failed requests:      0
15 Write errors:         0
16 Total transferred:    24822504 bytes
17 HTML transferred:     24574536 bytes
18 Requests per second:  12.54 [#/sec] (mean)
19 Time per request:     797.213 [ms] (mean)
20 Time per request:     79.721 [ms] (mean, across all requests)
21 Transfer rate:        603.31 [Kbytes/sec] received
22
23 Connection Times (ms)
24      min  mean[+/-sd] median   max
25 Connect:    48    65   14.3     61   145
26 Processing: 284   436  288.9    376  2957
27 Waiting:    119   199  281.5    151  2660
28 Total:      333   500  287.8    439  3007
29
30 Percentage of the requests served within a certain time (ms)
31  50%    439
32  66%    458
33  75%    477
34  80%    493
35  90%    563
36  95%    666
37  98%   1420
38  99%   2142
39 100%   3007 (longest request)
```

Listing 2.2 przedstawia wynik testów czasowych. Najważniejszą informacją z punktu widzenia optymalizacji jest ilość żądań na sekundę, która w tym wypadku wynosi 12.54. Narzędzie **ab** pozwala również zdiagnozować potencjalne błędy aplikacji pod wpływem zbyt dużego ruchu. Pola takie jak **Failed requests** oraz **Write errors** ułatwiają określenie prawidłowości wykonywania żądań. W powyższym przykładzie, wartości są akceptowalne (średni czas żądania to 0.5 sekundy), co najważniejsze nie występują błędy na poziomie serwera WWW i z dużym prawdopodobieństwem również na poziomie aplikacji. Oczywiście zauważalny jest spadek wydajności, w porównaniu z pierwszym testem, co prawda wartości średnie są zbliżone, jednak widać większe rozbieżności między wartościami minimalnymi a maksymalnymi. Najdłuższe zapytanie zajęło ponad 3 sekundy.

W dokumentacji aplikacji **ab**, można znaleźć informację, że niektóre serwery mogą blokować wysyłane przez niego nagłówki HTTP. W tym celu można wykorzystać przełącznik umożliwiający podanie się za inną przeglądarkę. Np. chcąc zasymulować odwiedziny przy użyciu przeglądarki Chrome

URL	Średnia	Min	Max	Błąd (%)	Req/Min
..14543704,wiadomosc.html	2299	415	94329	2.2	43.7
...1028235,wiadomosci.html	2485	335	94434	2.8	43.7
...14545325,wiadomosc.html	2023	394	94285	1.8	43.7
<b>Łącznie</b>	<b>2269</b>	<b>335</b>	<b>94434</b>	<b>2.27</b>	<b>131.1</b>

Tab. 2.1: Tabela z rezultatem działania aplikacji JMeter

należy wykonać poniższą komendę. `ab -n 100 -c 5 -H "Mozilla/5.0 (Windows; U; Windows NT 5.1; en-US) AppleWebKit/534.2 (KHTML, like Gecko) Chrome/6.0.447.0 Safari/534.2" http://www.example.com.`

Pomimo wielu zalet wynikających z korzystania z narzędzia `ab`, istnieje jedna zasadnicza wada. Aplikacja nie daje możliwości przetestowania pewnego scenariusza lub jest to bardzo niewygodne, a przypadku aplikacji z wykorzystaniem technologii *JavaScript* i *AJAX* wręcz niewykonalne.

Dlatego też na przełomie lat wyspecjalizowały się bardziej zaawansowane narzędzia przeznaczone do prześledzenia logicznej kolejności działań wykonywanych na stronie (pewnego przypadku użycia) i wykonywanie analiz właśnie w ramach logicznego zbioru akcji. W ten sposób możliwe jest przeprowadzenie tzw. *testów funkcjonalnych*, a także sprawdzenie w jakim stopniu są one wrażliwe na zwiększony ruch sieciowy.

Jednym z przykładów takiego narzędzia o naprawdę olbrzymich możliwościach jest *Apache JMeter*. Z jego pomocą możliwe jest nagranie pewnego ciągu akcji wykonanych przy pomocy przeglądarki, a następnie przeprowadzenie ciągu analiz i testów na tak wyodrębnionym zbiorze. Tabela 2.1 pokazuje wynik działania przykładowego scenariusza polegającego na symulowaniu wejścia na stronę `http://www.wp.pl`, a następnie kliknięciu jednego z linków wiadomości, po czym kliknięciu na kolejny link z dostępnych na bieżącej stronie. Ostatnim elementem łańcucha akcji jest wysłanie komentarza do artykułu. JMeter podobnie jak `ab` umożliwia wykonywanie równoległych połączeń użytkowników określanych jako wątki (*threads*).

W zdefiniowanym przypadku użycia 5 użytkowników wykonuje jednocześnie tą logikę 500 razy. Na zakończenie dane możliwe są do wyeksportowania do formatu *CSV* lub wyświetlone bezpośrednio na ekranie. JMeter jest narzędziem bardzo rozbudowanym, przez co idealnie nadaje się do testowania zaawansowanych scenariuszy zarówno pod kątem poprawności działania, jak również ogólnej wydajności.

## 2.2 Analiza wydajności frontendu

W poprzednim podrozdziale została omówiona tematyka testowania czasu działania zasobów serwowanych przez serwer WWW. Wykorzystując wspomniane wcześniej narzędzia można jednoznacznie ustalić czy aplikacja funkcjonuje w sposób prawidłowy, czy nie występują błędy w pracy serwera oraz jak dobrze oprogramowanie radzi sobie ze wzmożonym obciążeniem.

Wydawać by się mogło, że problem testowania wydajności aplikacji został jednoznacznie omówiony. Nic bardziej mylnego, jedynie w idealnym świecie, strona WWW składała by się wyłącznie z tekstu, a użytkownicy do przeglądania Internetu, korzystaliby wyłącznie z terminali.

Współczesny użytkownik internetu, do przeglądania jego zasobów wykorzystuje przeglądarkę internetową, zdolną do wyświetlania zarówno tekstu, jak i mediów wszelakiego typu. Dlatego też, dla większej precyzji, konieczne jest wyszczególnienie kluczowego komponentu oprogramowania, jakim jest *front-end* aplikacji.

W przypadku aplikacji internetowych, *front-end* to graficzny interfejs służący do komunikacji użytkownika ze stroną i prezentacji danych opracowanych przez zaplecze systemu (*back-end*) w sposób przystępny i zrozumiały.

Odpowiednia optymalizacja *front-endu* jest o tyle ważna, że jest to pierwsza technologia z jaką użytkownik ma kontakt w momencie korzystania z aplikacji webowej [2].

Frontend realizuje przetwarzanie i analizę wyniku odpowiedzi serwera, przy czym obowiązującym językiem komunikacji jest język HTML (*Hypertext Markup Language*). Od kilku lat w wyniku rozwoju trendu WEB 2.0, istotnym zabiegiem, projektowanych aplikacji, staje się przeniesienie części logiki na stronę przeglądarki (frontendu). Cienkie do tej pory aplikacje webowe (*thin client*), zaczynają - dzięki zdobyczą technologicznym takim jak *AJAX* realizują znacznie szerszy zakres funkcjonalny niż dotychczas. Oznacza to, że przetwarzanie danych może mieć miejsce wykorzystując przeglądarkę internetową i język JavaScript. Dlatego też kolejnym istotnym elementem analizy wydajnościowej staje się analiza *frontendu*.

Na podstawie badań firmy Juniper, stwierdzono, że średni czas, jaki użytkownik jest w stanie poczekać na załadowanie strony to zaledwie 4 sekundy. Nie warto więc tracić potencjalnych klientów strony, tylko i wyłącznie z powodu braku optymalizacji po stronie przeglądarki.

Wśród istniejących na rynku rozwiązań służących do analizy po stronie klienta, najpopularniejszymi są te, które są wbudowane bezpośrednio w interfejs przeglądarki internetowej (jest to przecież najbardziej intuicyjne podejście). Jednym z pierwszych rozwiązań tego typu była wtyczka **Firebug**

napisana dla przeglądarki Firefox. Jest to obecnie najbardziej zaawansowane narzędzie tego typu, rywalizujące jednocześnie z natywnymi dodatkami deweloperskimi dla przeglądarki Chrome.

Interfejs Firebuga pozwala na szczegółową inspekcję kodu HTML wraz z możliwością dynamicznej operacji na węzłach DOM dokumentu HTML. Nie mniej ważnymi narzędziami są: możliwość wykonywania i debugowania kodu JavaScript na stronie, inspekcja związanego z dokumentem HTML obiektu DOM, edycja i rewizja kodu JavaScript oraz narzędzie do monitorowania ruchu sieciowego wykonywanego przez aplikację.

Ten ostatni moduł pełni podobną rolę do narzędzia **ab**, jednak wyświetla wszystkie zasoby, które mają bezpośrednie powiązanie z bieżącym dokumentem HTML. Omawiane wcześniej narzędzia pokazywały jedynie czas renderowania dokumentu HTML, jednak należy mieć na uwadze, że strona internetowa składa się z wielu różnych zasobów, wśród których nie sposób pominąć: grafik, arkuszy CSS, kodu JavaScript, apletów Java czy obiektów Adobe Flash.

Każda strona może zawierać zróżnicowaną ilość takich zasobów, dlatego na łączny czas ładowania strony składa się zarówno czas oczekiwania na dokument HTML, jak również czas konieczny na pobranie każdego z powiązanych z nim zasobów.

Nawiązując do [3] istnieje zasada, która mówi, że tylko 10-20% czasu odpowiedzi jest spędzane na oczekiwanie dokumentu HTML, natomiast pozostałe 80-90% to czas na pobieranie pozostałych zasobów i ładowanie zawartości DOM.

Rysunek 2.3 jest najlepszym przykładem tej zasady. Strona wykonuje 32 zapytania do serwera, z czego tylko jedno to żądanie dokumentu HTML. Pobranie tego dokumentu zajęło około 400ms, tymczasem łączny czas wczytywania strony wyniósł **3.21 sekundy**. Oznacza to, że generowanie dokumentu zajęło zaledwie **12%** łącznego czasu oczekiwania.

Na podstawie danych z Firebuga łatwo stwierdzić pewne nieprawidłowości, bowiem w ramach strony wczytywane są 2 stosunkowo duże (1,1 MB) dokumenty graficzne, które prawdopodobnie nie zostały wymiarowane do odpowiednich rozmiarów. **Firebug** stanowi, więc cenne narzędzie przy diagnostyce *frontendu* strony. Narzędzie sstaje się jeszcze przydatniejszy przy pojawianiu się elementów dynamicznych JavaScript, ponieważ pozwala śledzić zarówno aktualnie wykonywany kod, jak również nasłuchiwać zapytań asynchronicznych wykonywanych przez AJAX. Narzędzie to może być również przydatne podczas śledzenia zmian dokonywanych w dokumencie HTML, za pomocą narzędzia inspekcji, umożliwia bowiem zbadanie każdego węzła.

## Rozdział 2. Analiza jakościowa i wydajnościowa aplikacji

URL	Status	Rozmiar	Oś czasu
GET www.ftims.p.lodz.p	200 OK	30.8 KB	392ms
GET ufo.js	304 Not Modified	11.1 KB	181ms
GET top.gif	304 Not Modified	37 B	60ms
GET switch_minus.gif	304 Not Modified	155 B	77ms
GET swf.swf?path=http:	304 Not Modified	313.1 KB	68ms
GET styles.php	200 OK	99.3 KB	456ms
GET styles.php	200 OK	13 KB	322ms
GET square.jpg	304 Not Modified	13.7 KB	316ms
GET shadow.png	304 Not Modified	4.8 KB	204ms
GET plakat.jpg	200 OK	1.1 MB	2.08s
GET pdf.gif	304 Not Modified	897 B	59ms
GET overlib_cssstyle.js	304 Not Modified	8.6 KB	137ms
GET overlib.js	304 Not Modified	48.1 KB	129ms
GET link.png	304 Not Modified	552 B	214ms
GET javascript-static.js	304 Not Modified	18.2 KB	114ms
GET javascript-mod.php	200 OK	34 B	264ms
GET headermain.png	304 Not Modified	115.8 KB	54ms
GET header_separator.p	304 Not Modified	308 B	247ms
GET header_separator.g	304 Not Modified	900 B	196ms
GET forumolddiscuss_s	304 Not Modified	237 B	266ms
GET footer.jpg	304 Not Modified	26.3 KB	316ms
GET f2.jpg	200 OK	1.6 KB	188ms
GET f2.jpg	200 OK	1.7 KB	253ms
GET dropdown.js	304 Not Modified	2.5 KB	191ms
GET dotted.gif	304 Not Modified	49 B	111ms
GET cookies.js	304 Not Modified	2.4 KB	168ms
GET content_separator.i	304 Not Modified	904 B	256ms
GET bottom.gif	304 Not Modified	35 B	51ms
GET borders.png	304 Not Modified	230 B	163ms
GET belge_box.png	304 Not Modified	3.6 KB	141ms
GET baner_Lodz_duzy.gl	200 OK	1.1 MB	2.1s
GET baner_FTIMS.JPG	200 OK	36.1 KB	382ms
32 requests			2.9 MB (572.3 KB z bufora podręcznego) 3.21s (onload: 3.21s)

Rys. 2.3: Analiza ruchu sieciowego na stronie <http://ftims.p.lodz.pl>



## 2.3 Analiza wydajności bazy danych

Bardzo często ogólna wydajność aplikacji uzależniona jest od szybkości operacji odczytu / zapisu bazy danych. W pewnym momencie twórcy aplikacji zaczynają oczekiwać od niej większej wydajności. Należy jednak zadać pytanie co należy tak naprawdę optymalizować? Konkretne zapytanie? Schemat bazy? Czy może sprzęt na którym baza danych pracuje? Jedynym sposobem na znalezienie jednoznacznej odpowiedzi, jest zmierzenie pracy wykonywanej przez bazę i sprawdzenie wydajności pod wpływem różnych czynników [1].

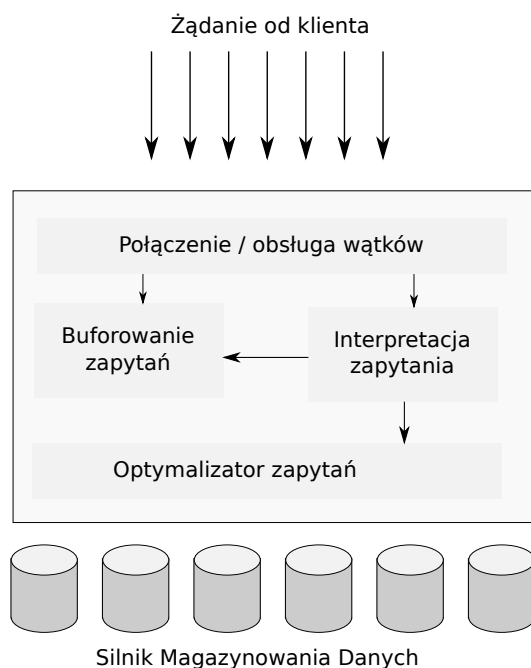
Bazy danych ewoluowały na przestrzeni kilkunastu lat, początkowo były to po prostu pliki o określonej strukturze, jednak wraz ze wzrostem wymagań zaczęto stosować równoległy dostęp do danych, a także dbając o spójność danych, zaimplementowano system transakcji.

Obecnie widać specjalizacji baz danych do konkretnych zastosowań, pomimo dominacji na rynku baz relacyjnych opartych o standard *SQL 93*, zaczęto również torować drogę nowym rozwiązaniom takim jak *NoSQL* czyli bazy danych o nieuporządkowanej strukturze, pozbawionych ściśle zdefiniowanych schematów, zyskując większą elastyczność. Dodatkowo w określonych zastosowaniach tego typu bazy danych okazują się dużo szybsze niż bazy relacyjne. Szybkość ta jest tym większa, im większa jest przechowywana kolekcja danych. Są to, więc bazy wysoce skalowalne, choć mniej spójne niż standardowe.

Wśród systemów bazodanowych wykorzystywanych w aplikacji e-commerce, bardzo dużą popularnością cieszy się oprogramowanie MySQL. Pomimo dużo mniejszych możliwości niż np. komercyjne rozwiązania Oracle czy MS SQL Server, omawiana baza danych jest elastyczna i łatwo zaadaptować ją do własnych potrzeb. Od czasu wprowadzenia zgodności ze standardem ACID, MySQL zaczął być szeroko wykorzystywany w e-biznesie.

Rysunek 2.4 przedstawia, jak wygląda architektura systemu baz danych MySQL, z punktu widzenia funkcjonalny komponentów [2][str. 26]. Pierwsza warstwa zawiera usługi, które wbrew pozorom nie są unikalne tylko dla omawianego oprogramowania. Są to usługi charakterystyczne dla większości narzędzi w architekturze sieciowej. Wyróżniono więc obsługę połączenia, uwierzytelnianie itd. Druga warstwa wprowadza wiele zmian i komponentów specyficznych dla MySQL'a. Wyróżniamy więc komponenty odpowiedzialne za parsowanie zapytań, a także ich optymalizację, buforowanie oraz kod odpowiedzialny za implementacje wbudowanych funkcji (np. data, czas, funkcje matematyczne i kryptograficzne). Każda funkcjonalność oferowana przez którykolwiek z silników składowania danych (*storage engine*), ma tu swoje miejsce (np. procedury użytkownika, wyzwalacze oraz widoki).

Trzecia warstwa wyróżnia wszystkie silniki składowania, które odpowiedzialne są bezpośrednio za przechowywanie i pobieranie danych w MySQL'u.



Rys. 2.4: Schemat architektury MySQL

Każdy z silników, ma różne zastosowania (podobnie jak różne systemy plików w systemach operacyjnych). Komunikacja z każdym z nich odbywa się wykorzystując wewnętrzne API. Wspomniany interfejs ukrywa różnice w specyfice każdego z mechanizmów, przez co zapytania są bardziej abstrakcyjne i prostsze w użyciu dla użytkownika końcowego. Podobnie jak w wypadku serwerów WWW, między mechanizmem składowania a serwerem MySQL występuje zależność: żądanie - odpowiedź. Oznacza to, że serwer wysyła żądanie i oczekuje na odpowiedź.

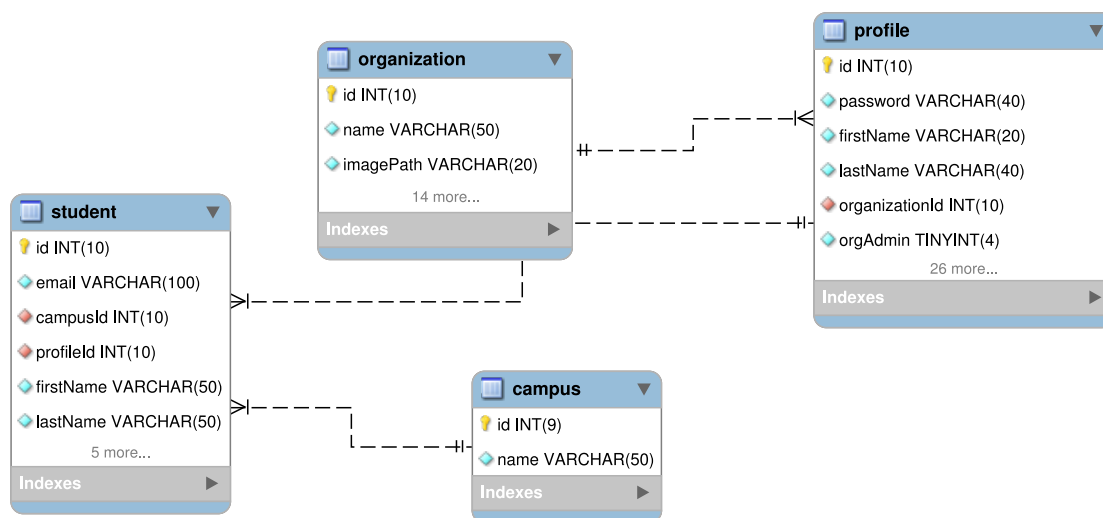
Najlepszą strategią optymalizacji, jest szukanie miejsc aplikacji, które działają najwolniej. Utworzono następujący schemat bazy danych (rys. 2.5). Na podstawie rysunku widać zależność studenta, który przynależy do 1 profilu (nauczyciela), który z kolei należy do organizacji. Podobnie student może należeć do jednego z istniejących kampusów.

Jak zachowa się baza danych przy próbie stworzenia alfabetycznego indeksu studentów, których nazwiska zaczynają się na daną literę (listing 2.3)?

#### Listing 2.3: Zapytanie do wyświetlenia menu książki adresowej uczniów

```
1 SELECT left(s.lastName, 1) letter, COUNT(s.id) students
2 FROM student s
3 WHERE s.lastName > 'A'
4 GROUP BY letter
```

Rezultat działania zapytania SQL widoczny jest na listingu 2.4. Jak wiadać, zapytanie wykonuje się w czasie 40 milisekund. Jest to dosyć krótko,



Rys. 2.5: Schemat testowej bazy danych

ale wynika to z głównie z rozmiarów kolekcji danych (20244 studentów). Wraz z rozrastaniem się tej tabeli, czas potrzebny na wykonanie tego zapytania będzie się systematycznie powiększał. Dzieje się tak, ponieważ tabela studentów nie posiada żadnego indeksu.

Listing 2.4: Wynik zapytania z listingu 2.3

1	+	-----	+	-----	+
2		letter		students	
3	+	-----	+	-----	+
4		A		999	
5		B		1749	
6		C		814	
7		D		670	
8		E		385	
9		F		918	
10		G		1599	
11		H		808	
12		I		245	
13		J		194	
14		K		1695	
15		L		1156	
16		M		1411	
17		N		492	
18		O		240	
19		P		686	
20		Q		14	
21		R		1259	
22		S		2662	
23		T		499	
24		U		57	
25		V		244	
26		W		698	
27		Y		316	

```

28 | Z          |          400 |
29 +-----+-----+
30 25 rows in set (0.04 sec)

```

## Czym jest indeks?

*Indeks* jest strukturą danych przeznaczoną do pomocy systemowi baz danych w efektywnym pobieraniu informacji z tabel. Są one często wymagane dla zapewnienia dobrej wydajności. Indeksy są szczególnie ważne w momencie kiedy baza danych się rozrasta, ponieważ ilość elementów do przeszukiwania wierszowego ulega zwielokrotnieniu.

Podczas zwykłego wyszukiwania wartości w bazie danych, program musi przeszukać każdą kolumnę, każdego wiersza w poszukiwaniu określonej wartości. W wypadku indeksów sprawa jest uproszczona ponieważ dysponujemy pewnym podzbiorem wartości np. z danej kolumny lub wyrażenia. W ten sposób silnik bazy danych wie, że dana wartość znaleziona w określonym indeksie powiązana jest z określonym rekordem, więc odpowiedź jest błyskawiczna.

Oczywiście wykorzystywanie indeksów wiąże się również z pewną zajętością danych, ponieważ oprócz danych w tabeli, trzeba dodatkowo przechowywać dane indeksów. Przeliczając jedna zyski do strat, większość przemawia jednak za stosowaniem indeksów.

### 2.3.1 Optymalizacja table przy użyciu indeksów

Optymalizacja zapytań przy użyciu indeksów jest stosunkowo prosta i polega na stworzeniu indeksu, który najlepiej pasuje do wyszukiwanej zawartości. W naszym wypadku potrzeba indeksu przechowującego pierwszą literę nazwiska. Z drugiej strony sortowanie po nazwisku lub nawet wyszukiwanie po nim jest dość częstą operacją dlatego warto stworzyć kompletny indeks dla pola `lastName` tabeli `student` (2.5). Na podstawie wyniku uzyskanego w listingu 2.6, widać znaczną poprawę czasu wykonywania. Wszystkie zapytania wykorzystujące kolumnę `lastName`, powinny wykonywać się zdecydowanie szybciej.

Listing 2.5: Utworzenie indeksu na polu nazwiska dla tabeli `student`

```
1 ALTER TABLE 'student' ADD INDEX 'lastName_idx' (('lastName');
```

Listing 2.6: Wynik zapytania z listingu 2.3 po optymalizacji indeksu

```

1 +-----+-----+
2 | letter | students |
3 +-----+-----+
4 | A      |          999 |
5 | B      |         1749 |
6 | C      |          814 |
7 | D      |          670 |

```

```
 8 | E      |      385 |
 9 | F      |      918 |
10 | G      |     1599 |
11 | H      |      808 |
12 | I      |      245 |
13 | J      |      194 |
14 | K      |     1695 |
15 | L      |     1156 |
16 | M      |     1411 |
17 | N      |      492 |
18 | O      |      240 |
19 | P      |      686 |
20 | Q      |        14 |
21 | R      |     1259 |
22 | S      |     2662 |
23 | T      |      499 |
24 | U      |        57 |
25 | V      |      244 |
26 | W      |      698 |
27 | Y      |      316 |
28 | Z      |      400 |
29 +-----+
30 25 rows in set (0.01 sec)
```

### Kolejność złączeń

Złączenia (ang. *joins*), są operacją spajającą ze sobą dwie lub więcej tabel. Dopuszczalne są łączenia 2 różnych tabel lub tej samej (*self-joins*). Używanie złączeń wynika w dużej mierze z normalizacji bazy danych, a co za tym idzie usunięcia nadmiarowości z rekordów tabel. Takie podejście w dużej mierze poprawia spójność danych, ale również może pogorszyć w dużej mierze wydajność zapytań.

Standard ANSI wyróżnia 4 rodzaje złączeń: *INNER*, *OUTER*, *LEFT*, *RIGHT*. W specjalnych okolicznościach tabela może być również połączona sama ze sobą. Zasadnicza różnica w specyfice polega na kryterium łączenia - w wypadku *INNER JOINÓW* wymagane jest istnienie odpowiadających sobie krotek po dwóch stronach relacji, natomiast *OUTER JOIN* wymaga spełnienia tego kryterium przynajmniej na jednej ze strony relacji - odpowiednio lewej lub prawej.

Różnice między złączeniami przekładają się również na wydajność działania, najpopularniejsze ze złączeń *inner join'y* są najszybsze, podczas, gdy pozostałe przeznaczone są do bardziej specyficznych zastosowań.

Kolejny z przeprowadzanych testów, będzie polegał na złączeniu ze sobą studentów oraz profili, innymi słowy wyświetleniu wszystkich studentów przynależnych do któregoś z profili. Wynik działania poszczególnych testów (listing 2.7) został zestawiony w tabeli 2.2. Wyniki zapytań stanowią potwierdzenie ogólnie przyjętych zasad optymalizacji zapytań, *inner join* okazał się najszybszym ze złączeń, jednocześnie widać, że kolejność wykonywania złączeń również ma wpływ na wydajność. Zazwyczaj powinno się

Typ	Czas wykonywania [sec]	Ilość rekordów
inner join student	0.2557	20244
left join student	<b>5.2125 / 0.5635</b>	20463
right join student	0.2634	20244
inner join profile	0.2726	20244
left join profile	0.3272	20463
right join profile	<b>5.2236 / 0.5704</b>	20463

Tab. 2.2: Tabela z rezultatem działania poszczególnych zapytań

zaczynać od tabeli, która posiada mniej rekordów (w tym wypadku tabela *profile*). Ogromna różnica między czasem wykonywania *left joina* oraz *right joina* w odwrotnej kolejności łączenia wynika z braku indeksu dla pola *profileId* tabeli *student*. Po dodaniu indeksów widać 10-krotną poprawę szybkości wykonywania tego zapytania.

Przedstawione dotychczas przypadki były stosunkowo proste do naprawy, często jednak zapytania są dużo bardziej rozbudowane i ciężkie do szybkiej dekompozycji. Warto wtedy skorzystać z udostępnianego przez MySQL narzędzia *EXPLAIN*. Oferuje on pomoc w zakresie dekompozycji bardziej skomplikowanych zapytań. Narzędzie to pokazuje m.in. wykorzystane klucze dla łączeń, liczebność łączonych tabel, proponowane usprawnienia indeksów.

## Listing 2.7: Kilka możliwych do wykorzystania zapytań

```

1 SELECT SQL_NO_CACHE * FROM profile p inner join student s on s.
   profileId = p.id
2 SELECT SQL_NO_CACHE * FROM profile p left join student s on s.
   profileId = p.id
3 SELECT SQL_NO_CACHE * FROM profile p right join student s on s.
   profileId = p.id
4
5 SELECT SQL_NO_CACHE * FROM student s inner join profile p on s.
   profileId = p.id
6 SELECT SQL_NO_CACHE * FROM student s left join profile p on s.
   profileId = p.id
7 SELECT SQL_NO_CACHE * FROM student s right join profile p on s.
   profileId = p.id

```

Jednym z zapytań, które może stanowić potencjalny problem w analizie jest zapytanie zaczerpnięte z istniejącej aplikacji opartej o przedstawiony na rysunku 2.5 schemat bazy danych. Przedstawione na listingu 2.8 zapytanie służy do pokazania wartości sum w poszczególnych aktywnościach takich jak programy studenckie, praktyki, ilość studentów itp. Zapytanie to jest wykonywane w kontekście określonego roku szkolnego, a także konkretnej organizacji - wyświetla wartości dla podległych kampusów. Omawiane za-

pytanie jest prawdopodobnie jednym z trudniejszych, na jakie można kiedykolwiek trafić. Głównym problemem takich zapytań jest istnienie wielu powiązanych z nim *zapytań skorelowanych*. W celu analitycznej analizy tego zapytania użyto narzędzie EXPLAIN. Składnia tego narzędzia jest prosta i polega na dodaniu tej dyrektywy do wcześniejszego zapytania. Wynik takiego zapytania zwróci ilość wierszy równą ilości podzapytań wykonywanych w trakcie przetwarzania (??).

Listing 2.8: Bardziej rozbudowane zapytanie SQL

```
1 SELECT
2 o.id, o.name,
3 (SELECT group_concat(oc.campusId SEPARATOR ' ' ) FROM
   organizationCampus oc WHERE oc.organizationId = o.id) as
   campusIds,
4 (SELECT count(s5.id)
5 FROM profile p5
6 INNER JOIN student s5 on s5.profileId = p5.id
7 WHERE p5.organizationId = o.id) as students,
8
9 IFNULL((SELECT count(DISTINCT s3.id)
10 FROM profile p3
11 INNER JOIN student s3 on p3.id = s3.profileId
12 INNER JOIN studentIntensiveProgram sip on sip.studentId = s3.id
13 WHERE p3.organizationId = o.id
14 GROUP BY p3.organizationId
15 ),0) as programs,
16
17 IFNULL((SELECT round(SUM(classes*1 + 1on1*3 + shabbaton*5 +
   socialEvents*0.5 + shabbosMeals*2))
18 FROM
19 student s2
20 INNER JOIN profile p2 on p2.id = s2.profileId
21 INNER JOIN 'reportStudentAttendance' AS 'r1' ON r1.studentId = s2.
   id
22 INNER JOIN 'report' AS 'ra' ON ra.id = r1.reportId and ((ra.month
   >= 9 and ra.year = 2011)
23 or (ra.month <=8 and ra.year = 2012))
24 WHERE p2.organizationId = o.id
25 GROUP BY p2.organizationId ), 0) as score,
26
27 (SELECT ifnull(sum(datediff("2012-08-30", sy.startDate) BETWEEN 30
   and 90
28 and sy.startDate between "2011-09-01" and "2012-08-30"),0)
29 FROM profile p4
30 INNER JOIN student s4 on s4.profileId = p4.id
31 INNER JOIN studentYeshiva sy on sy.studentId = s4.id
32 WHERE p4.organizationId = o.id) as yeshiva_1_3,
33
34 (SELECT ifnull(sum(datediff("2012-08-30", sy.startDate) BETWEEN 91
   and 180
35 and sy.startDate between "2011-09-01" and "2012-08-30"),0)
36 FROM profile p4
37 INNER JOIN student s4 on s4.profileId = p4.id
38 INNER JOIN studentYeshiva sy on sy.studentId = s4.id
39 WHERE p4.organizationId = o.id) as yeshiva_4_6,
40
41 (SELECT ifnull(sum(datediff("2012-08-30", sy.startDate) > 181
42 and sy.startDate > 2000),0)
```

```
43 FROM profile p4
44 INNER JOIN student s4 on s4.profileId = p4.id
45 INNER JOIN studentYeshiva sy on sy.studentId = s4.id
46 WHERE p4.organizationId = o.id) as yeshiva_6,
47
48 IFNULL((SELECT sum(case
49 when s1.beganSo = "spring/2011" then "2011-01-01"
50 when s1.beganSo = "summer/2011" then "2011-05-01"
51 when s1.beganSo = "fall/2011" then "2011-09-01"
52 else "2012-08-30"
53 end >= "2011-09-01" and right(s1.beganSo, 4) >= 2011 and s1.beganSo
    != "before")
54 FROM student s1
55 INNER JOIN profile p1 on s1.profileId = p1.id
56 WHERE p1.organizationId = o.id
57 GROUP BY o.id),0) AS 'so'
58 FROM organization o
59 GROUP BY o.id
60 HAVING score > 0
61 ORDER BY o.name
```



id	select_type	table	type	possible_keys	key
1	PRIMARY	o	ALL		
9	DEPENDENT SUBQUERY	p1	ref	PRIMARY,organizationId	organizationId
9	DEPENDENT SUBQUERY	s1	ref	profile_idx	profile_idx
8	DEPENDENT SUBQUERY	p4	ref	PRIMARY,organizationId	organizationId
8	DEPENDENT SUBQUERY	s4	ref	PRIMARY,profile_idx	profile_idx
8	DEPENDENT SUBQUERY	sy	ref	studentIdx	studentIdx
7	DEPENDENT SUBQUERY	p4	ref	PRIMARY,organizationId	organizationId
7	DEPENDENT SUBQUERY	s4	ref	PRIMARY,profile_idx	profile_idx
7	DEPENDENT SUBQUERY	sy	ref	studentIdx	studentIdx
6	DEPENDENT SUBQUERY	p4	ref	PRIMARY,organizationId	organizationId
6	DEPENDENT SUBQUERY	s4	ref	PRIMARY,profile_idx	profile_idx
6	DEPENDENT SUBQUERY	sy	ref	studentIdx	studentIdx
5	DEPENDENT SUBQUERY	p2	ref	PRIMARY,organizationId	organizationId
5	DEPENDENT SUBQUERY	s2	ref	PRIMARY,profile_idx	profile_idx
5	DEPENDENT SUBQUERY	r1	ref	student_idx,report_idx	student_idx
5	DEPENDENT SUBQUERY	ra	eq_ref	PRIMARY,month_idx,year_idx	PRIMARY
4	DEPENDENT SUBQUERY	p3	ref	PRIMARY,organizationId	organizationId
4	DEPENDENT SUBQUERY	s3	ref	PRIMARY,profile_idx	profile_idx
4	DEPENDENT SUBQUERY	sip	ref	student_idx	student_idx
3	DEPENDENT SUBQUERY	p5	ref	PRIMARY,organizationId	organizationId
3	DEPENDENT SUBQUERY	s5	ref	profile_idx	profile_idx
2	DEPENDENT SUBQUERY	oc	index		PRIMARY

Tab. 2.3: Wynik działania polecenia EXPLAIN

## Rozdział 3

# Wymagania i budowa aplikacji

Aplikacja powinna spełniać wszystkie oczekiwane cele biznesowe, jak również zapewniać prawidłowe funkcjonowanie bez względu na aktualnie panujące warunki. Aplikacja demonstracyjna będzie księgarnią internetową oferującą wiele typowych funkcji, charakterystycznych dla tej branży.

Głównym celem tworzonego oprogramowania jest odpowiednie wyważenie pracy wykonanej po stronie serwera, jak również po stronie przeglądarki. Dlatego też duży nacisk pracy został położony na utworzenie usług udostępniających interfejs dostępowy do bazy danych. W ten sposób wykorzystując architekturę REST kod JavaScript może dokonywać modyfikacji modelu, przy najmniejszym udziale serwera WWW.

### 3.1 Zakres funkcjonalny

Projektowana aplikacja stanowi wirtualną księgarnię, podobnie jak w wypadku jej realnego odpowiednika, zapewnia możliwość przeglądania zasobów, podzielonych na kategorie lub oznaczonych określonymi słowami kluczowymi. W odróżnieniu od prawdziwej księgarni, w internetowych aplikacjach konieczne jest zdefiniowanie pewnej tożsamości, która będzie później podstawą do zakupu książki, lub sprawdzenia stanu zamówienia.

Dlatego też istotne jest stworzenie spójnego systemu uwierzytelniania, zapewniającego zarówno bezpieczeństwo, jak również łatwość ewentualnego przypomnienia zapomnianego hasła.

Ponieważ oferta książkowa, ciągle się zmienia, na rynek trafiają nowe książki, a nawet tworzą się nowe gatunki, istotną rolę w tworzonej aplikacji stanowi współlistnienie zarówno części przeznaczonej dla zwykłego użytkownika, jak również części administracyjnej (*backend*). Sekcja ta powinna być również zabezpieczona przed potencjalnym włamaniem lub złośliwym atakiem ze strony intruzów.

Część Administracyjna	Sekcja użytkownika
Dodawanie książek	Rejestracja
Edycja książek	Edycja i podgląd konta
Usuwanie książek	Ocenianie książek
Moderacja wpisów	Komentowanie książek
Dodawanie kategorii	Przeglądanie kategorii
Dodawanie plików do książek	Pobieranie streszczeń książek
Nadzorowanie kont użytkowników	Wyszukiwanie książek w wyszukiwarce

Tab. 3.1: Tabela z wykazem funkcjonalności

Ponieważ obecnie w dobie ery WEB 2.0 strony oferują dużo większą interakcję, niż było to możliwe na początku ery Internetu. Dlatego też projektowana aplikacja daje również możliwość oceny i komentowania istniejących zbiorów. Jest to bardzo ważne z punktu widzenia biznesu i sprzedaży, ponieważ popularność i przychylne noty, z pewnością napędzą chętnych do kupna.

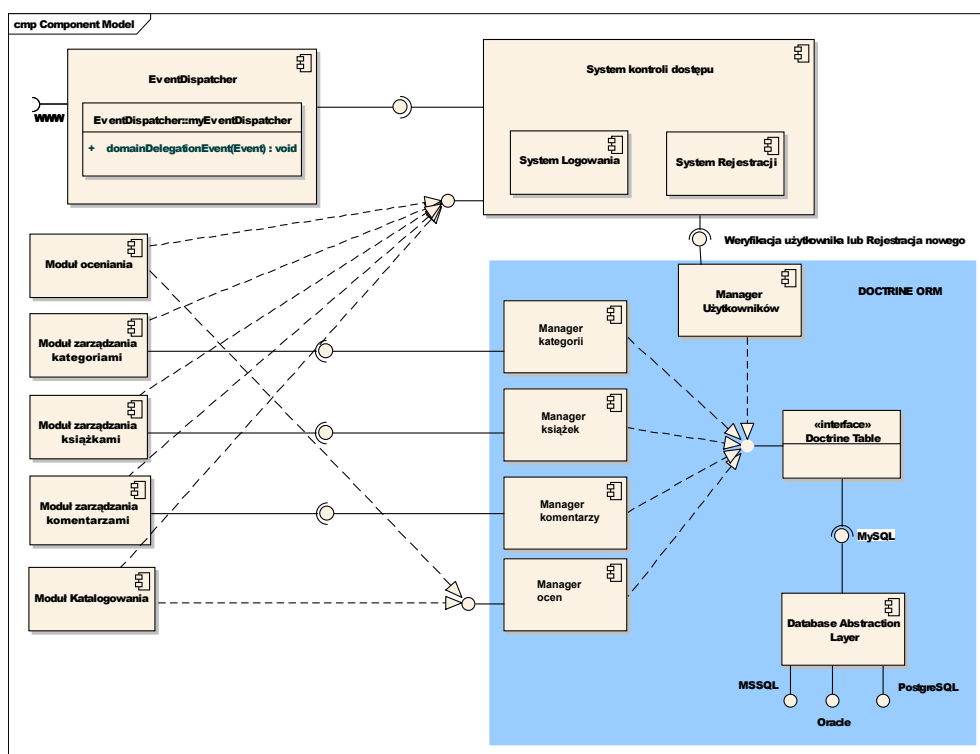
Wymagania podzielono na część administracyjną i część użytkownika, natomiast zbiór funkcjonalności przedstawiony jest w tabeli 3.1.

Diagram komponentów dla tworzonej aplikacji został uwzględniony na rysunku 3.1. Jak można wywnioskować z diagramu, aplikacja wykorzystuje mapowanie obiektowo relacyjne w celu ułatwienia komunikacji z bazą danych. Jest to bardzo wygodne rozwiązanie dające wiele swobody i elastyczności tworzonemu oprogramowaniu.

Aplikacja powinna w jak największym stopniu odciążać serwer WWW od niepotrzebnych zapytań, w tym celu, akcje usuwania, oceniania i komentowania książek będą realizowane przez kod JavaScript. Dodatkowo część walidacji danych będzie w pierwszej fazie wykonywana również przy użyciu frontendu, dopiero w wypadku wyłączenia wykonywania skryptów, wykonane zostanie żądanie do sprawdzenia przez serwer. Takie podejście powinno znacznie zredukować obciążenie np. podczas rejestracji użytkowników, ponieważ serwer wykonuje tylko minimum swojej pracy (*Lazy Loading*).

Należy jednak mieć na uwadze bezpieczeństwo oprogramowania, dlatego nie należy polegać w 100% na walidacji wykonywanej przez JavaScript, ponieważ użytkownik może bez problemu wyłączyć działanie skryptów. Dlatego też podczas tworzenia tego typu rozwiązań walidacja po stronie klienta powinna iść zawsze w parze z walidacją po stronie serwera.

Ze względu na konserwację i rozwój tworzonej aplikacji, jest ona zarządzana przy pomocy systemu kontroli wersji git. W ten sposób tworzony kod



Rys. 3.1: Diagram komponentów tworzonej aplikacji

jest pod ścisłą kontrolą, może być potem rozwijany przez więcej niż jednego programistę, eliminując konflikty podczas pracy na wspólnych zasobach.

Aplikacja wykorzystuje język HTML5, który staje się powoli standardem tworzenia kodu po stronie klienta. Zawiera bowiem wiele nowoczesnych mechanizmów i udogodnień w stosunku do poprzedniej wersji np. HTML4 oraz XHTML 1.0. Mając na uwadze prawidłowe funkcjonowanie w poszczególnych przeglądarkach i różną implementację nowego standardu, aplikacja wykorzystuje dodatek **HTML5 Boilerplate** oraz **Twitter Bootstrap** jako metodę na otrzymanie podobnych rezultatów w szerokiej gamie urządzeń (w tym urządzenia mobilne).

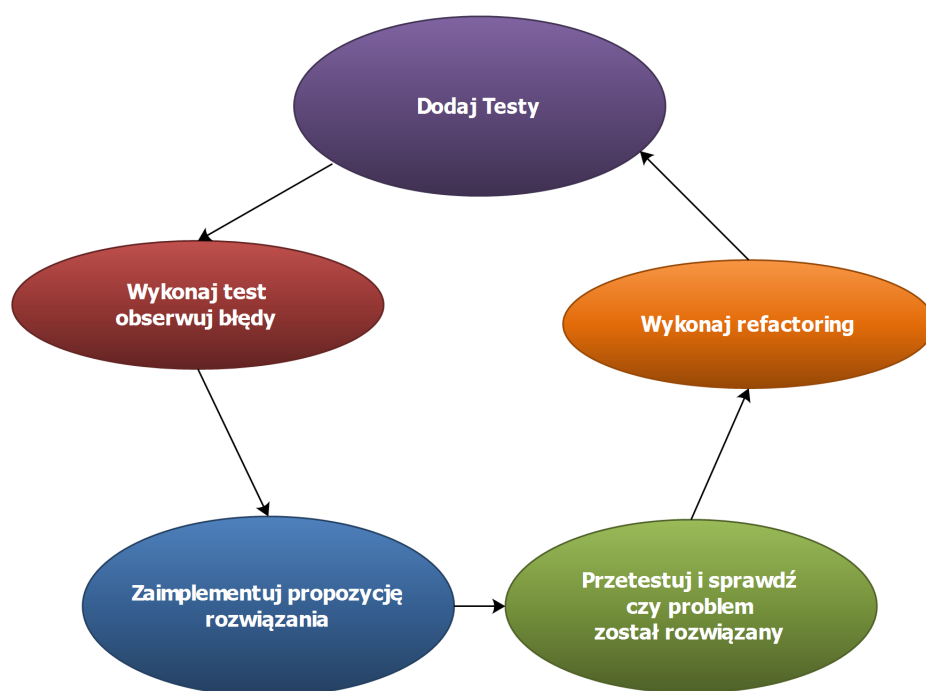
Projektowana aplikacja wykorzystuje technologię AJAX (ang. Asynchronous JavaScript and XML, asynchroniczny JavaScript i XML), dzięki czemu możliwa jest symulacja zachowania aplikacji desktopowej. Reakcja na akcje użytkownika następuje bez potrzeby przeładowania strony. Pozwala to na dużą oszczędność przepustowości i czasu koniecznego do przetworzenia całej strony. Dzięki zastosowaniu frameworka jQuery dla języka JavaScript, możliwe jest stosowanie bardziej przyjaznego interfejsu użytkownika, przy jednoczesnym zapewnieniu kompatybilności wstecz z istniejącymi wersjami przeglądarek internetowych.

### 3.2 Budowa aplikacji

Aplikacja została zbudowana przy pomocy frameworka symfony wykorzystując wzorzec projektowy MVC. W ten sposób możliwa jest separacja warstw logiki, modelu oraz widoku. Zapewnia to łatwość w utrzymaniu oprogramowania, przez co kod jest bardziej zrozumiały i mniej podatny na błędy.

Każdy element w systemie posiada szereg zależności, których nieprawidłowe działanie pociąga wiele daleko idących konsekwencji. Dlatego ważne jest, by kod aplikacji był spójny i podzielony na mniejsze jednostki budulcowe. Umożliwi to łatwiejsze wykrywanie problematycznych aspektów aplikacji, istotnych dla prawidłowego działania całego projektu. Ręczne testowanie całego projektu jest procesem żmudnym, trudnym, a niekiedy nawet niemożliwym. Dlatego istotne jest, by w rozbudowanych systemach z obszernym drzewem zależności, móc łatwo testować wprowadzane zmiany i przewidywać możliwe skutki uboczne.

W rozbudowanych aplikacjach zazwyczaj tylko osoba architekta oprogramowania ma całościowe pojęcie o tworzonej aplikacji. Dlatego też zwykły programista tworzący dany fragment aplikacji, może nie być świadom skutków ubocznych opracowywanej implementacji w stosunku do całego projektu. Ciągła konsultacja tworzonego kodu z architektem jest nie tylko nieekonomiczna, ale także może się okazać nieskuteczna. Dzieje się tak, bowiem



Rys. 3.2: Cykl wytwarzania w metodyce TDD

pewne detale działania aplikacji można omyłkowo pominąć lub zlekceważyć.

Rozwiązanie omawianego problemu polega na wykorzystaniu techniki tworzenia oprogramowania TDD (ang. *Test-driven Development*). Technika *programowania sterowanego testami* jest zaliczana do metodyk zwinnych *Agile*. Na rysunku 3.2, pokazano typowy cykl wytwarzania aplikacji. Każdy mechanizm wytwarzanego oprogramowania zakłada *pokrycie testami*, a same testy pisane są przed faktycznym kodem aplikacji.

Bardzo często aplikacja musi ulec modyfikacją polegającym na relokacji pewnych fragmentów kodu, przenosząc je w bardziej odpowiednie dla nich miejsce. Podczas takich zmian funkcjonalność kodu nie ulega zmianie, natomiast można narzucić realizację pewnych standardów, określanych mianem *konwencji kodowania*. Technika zmian określana jako *refactoring*, jest sposobem na zwiększenie czytelności kodu i polepszenie jego logicznej konstrukcji. Istnieje jednak ryzyko błędów wynikających z zależności pomiędzy modyfikowanymi fragmentami kodu a całą aplikacją. Stosując metodykę *TDD*, można uniknąć za każdym razem prawdopodobieństwa wystąpienia błędów wynikających z braku spójności.

### 3.3 Projekt bazy danych

Rysunek 3.3 przedstawia projekt bazy danych wykorzystanej przy okazji projektu. Diagram ERD przedstawia jednak bazę w podejściu relacyjnym, w późniejszej części pracy pokazana zostanie baza NoSQL, oraz proces denormalizacji, który musiał zostać dokonany.

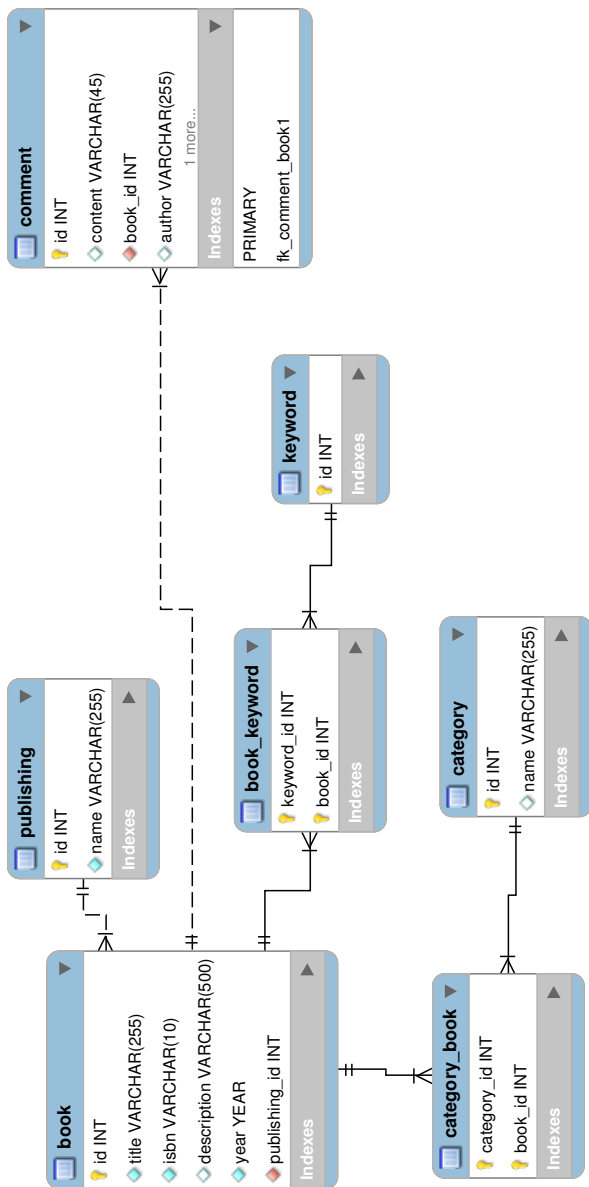
### 3.4 Organizacja kodu

Pewnym ułatwieniem w implementacji projektu będzie stworzenie mapy organizacji kodu, która umożliwi umiejscowienie elementów budulcowych projektu w kontekście całej aplikacji. Wyróżniono podział aplikacji, modułów oraz najbardziej atomowej części projektu, czyli akcji. Każdy z modułów wydziela charakterystyczną dla siebie dziedzinę przetwarzania, w większości wypadków ograniczaną przez docelową klasę modelu. W ten sposób moduł `school` dla aplikacji `frontend`.

Dzięki zastosowaniu technologii tworzenia prototypów modułów na podstawie istniejącego schematu bazy (*scaffolding*), większość modułów wchodzących w skład aplikacji `backend`, mogła zostać utworzona w sposób zautomatyzowany. W ten sposób implementacja modułów dla tej aplikacji została znormalizowana do interfejsu *CRUD*, a dzięki systemowi obsługi zdarzeń, można nadać dla każdego z modułów uniwersalne restrykcje. Budowa i działanie *generatorów*, wykorzystanych do osiągnięcia wspomnianego rozwiązania, zostało omówione w rozdziale ??.

Specyfika projektu zakłada konieczność komunikacji pomiędzy wieloma niezależnymi komponentami, które składają się na framework *Symfony*. Stwarza to konieczność utworzenia diagramu komponentów. Na podstawie rysunku 3.1 oraz przedstawionych w tym rozdziale przypadków użycia wiadać, że aplikacja wymaga zastosowania rozbudowanego systemu uprawnień *ACL* (*Access Control List*). Ponadto, konieczne jest przełożenie adresów URL, na odpowiednią logikę aplikacji. W tym celu wykorzystano komponent `EventDispatcher`, który dokona odpowiedniej interpretacji nadchodzących żądań.

Ponieważ projektowana aplikacja stawia wiele wymagań w stosunku do realizacji, konieczne będzie zastosowanie wzorca *Model-Widok-Kontroler* oraz narzędzi gwarantujących łatwy dostęp do bazy danych celem zapewnienia spójności implementacji.



Rys. 3.3: Projekt bazy danych



## Rozdział 4

# Architektura aplikacji

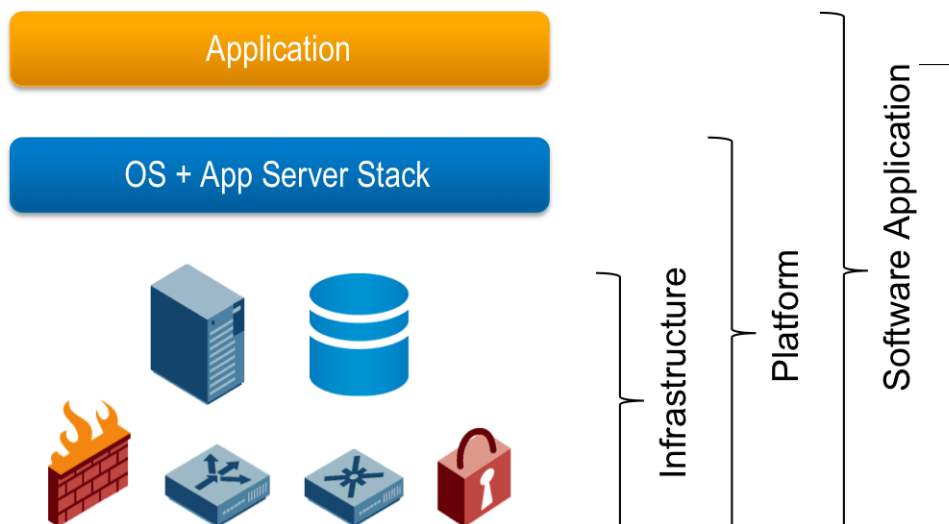
Istotnym celem pracy jest zapewnienie możliwie najlepszej w danej chwili wydajności. Twierdzenie to powinno być prawdziwe również wówczas, kiedy aplikacja znajduje się pod działaniem silnego ruchu sieciowego. Aby to zapewnić, konieczne jest wykorzystanie odpowiedniej architektury aplikacji.

Od kilku lat na rynku IT można zaobserwować duże zainteresowanie związane z chmurami obliczeniowymi. Jeszcze większe zamieszanie na rynku spowodowało udostępnienie przez Google oraz Microsoft ich produktów znanych jako *Google Application Engine* oraz *Windows Azure*. Są to usługi kwalifikowane jako *PaaS* czyli *Platform as a Service*. Oznacza to, że korzystając z ich produktów otrzymujemy kompletne środowisko uruchomieniowe aplikacji oraz zaplecze technologiczne umożliwiające uruchamianie aplikacji. W wypadku GAE możliwe jest programowanie aplikacji z wykorzystaniem udostępnionego przez usługę zbioru bibliotek, napisanych w trzech językach: Java, Python oraz Go. Następnie przy pomocy udostępnionych narzędzi możliwe jest wdrożenie aplikacji (*Deployment*).

Takie podejście do tworzenia aplikacji internetowych zyskało olbrzymią ilość zwolenników, ponieważ pozwoliło na całkowite przeniesienie ciężaru zarządzania skomplikowaną infrastrukturą na producentów rozwiązań *PaaS*. Innym ważnym powodem dla którego wielu ludzi zdecydowało się na wykorzystanie nowych usług, jest możliwość konfiguracji architektury do własnych potrzeb, a także (co było kluczowym powodem), do aktualnego obciążenia aplikacji.

### 4.1 Rodzaje chmur

Rysunek 4.1 pokazuje zasadnicze różnice między poszczególnymi rodzajami chmur. Chmura oferowana przez Google Application Engine, zapewnia wsparcie w zakresie architektury serwera, systemu plików, a także systemu bazodanowego. Dodatkowo GAE udostępnia również możliwość korzystania z usług w tle oraz serwerów mailowych.



Rys. 4.1: Zestawienie rodzajów chmury ze względu na udostępniane zasoby

Jak, więc wynika z zestawienia rodzajów chmur, oferuje ona dużo więcej ponad standardową infrastrukturę.

## 4.2 Nie relacyjne bazy danych NoSQL

Wraz z wykorzystaniem nie relacyjnych baz danych, zmienił się zupełnie pomysł na organizację danych. Siłą baz danych NoSQL jest brak ściśle określonej struktury, przez co schemat danych może się dowolnie zmieniać w trakcie rozwoju aplikacji. Innym ważnym powodem wykorzystania bazy Google Datastore oferowanej przez GAE jest szybkość i skalowalność. W odróżnieniu od zwykłego hostingu, baza danych w GAE może być rozproszona na dowolną ilość instancji.

## 4.3 Aplikacja zorientowana na usługi

Inną ważną cechą tworzonej aplikacji, jest wyszczególnienie usług (*webservicesów*), które pomogą w wykonywaniu operacji na modelu przy pomocy metod protokołu HTTP. Będą to więc popularne w dzisiejszych czasach usługi REST. Ze względu na specyfikę aplikacji i duży udział języka JavaScript w tworzonej logice, w większości wypadków usługi te będą zwracały w notacji JSON.

## Rozdział 5

# Optymalizacja aplikacji

Optymalizacja

## Rozdział 6

# Optymalizacja kodu klienta

Optymalizacja kodu klienta

## Rozdział 7

# Metody rozproszenia aplikacji i usług

Rozproszenie

# Podsumowanie

Podsumowanie

# Bibliografia

- [1]
- [2] A. Padilla and T. Hawkins. *Pro PHP Application Performance Tuning PHP Web Projects for Maximum Performance*. Apress.
- [3] S. Souders. *High Performance Web Sites*. O'Reily.

# Spis rysunków

2.1	Analiza czasu wykonywania strony <a href="http://ftims.edu.p.lodz.pl/">http://ftims.edu.p.lodz.pl/</a> wykonana w przeglądarce Google Chrome . . . . .	6
2.2	Cykl życia żądania . . . . .	8
2.3	Analiza ruchu sieciowego na stronie <a href="http://ftims.p.lodz.pl">http://ftims.p.lodz.pl</a> . .	14
2.4	Schemat architektury MySQL . . . . .	16
2.5	Schemat testowej bazy danych . . . . .	17
3.1	Diagram komponentów tworzonej aplikacji . . . . .	26
3.2	Cykl wytwarzania w metodyce TDD . . . . .	28
3.3	Projek bazy danych . . . . .	30
4.1	Zestawienie rodzajów chmury ze względu na udostępniane za- soby . . . . .	32



# Spis listingów

2.1	Analiza strony z wykorzystaniem narzędzia <b>ab</b> . . . . .	7
2.2	Test obciążenia czasowego . . . . .	9
2.3	Zapytanie do wyświetlenia menu książki adresowej uczniów .	16
2.4	Wynik zapytania z listingu 2.3 . . . . .	17
2.5	Utworzenie indeksu na polu nazwiska dla tabeli student . . .	18
2.6	Wynik zapytania z listingu 2.3 po optymalizacji indeksu . . .	18
2.7	Kilka możliwych do wykorzystania zapytań . . . . .	20
2.8	Bardziej rozbudowane zapytanie SQL . . . . .	21

# Płyta CD

Wraz z treścią pracy dyplomowej dołączono również płytę CD z kompletnym kodem źródłowym aplikacji. Dodatkowo kod można pobrać z poniższego repozytorium SVN: