

Simulace Škálování webové služby - Reaktivní vs Prediktivní přístup

Datum:

27.11.2024

Autoři:

Jakub Fukala - xfukal01

Adam Kozubek - xkozub09

1. Úvod

Smyslem této práce je vytvořit simulaci [10] škálování webové služby při proměnných podmínkách. Simulace bude porovnávat dva přístupy škálování (reaktivní a prediktivní). Na základě výsledků simulace vzorových modelů jsme stanovili obecné poznatky a demonstrovali za jakých okolností je lepší využít který přístup. Vytvořený simulační nástroj je ale určen primárně pro firmy na pomoc při rozhodování, kterou podobu škálování a s jakým nastavením je vhodné použít pro konkrétní službu.

Bylo nutné okrajově nastudovat fungování algoritmů Autoscalingu, zjistit reálné režie jednotlivých procesů a především vytvořit co nejpřesnější model [11] předpokládané vytíženosti a na základě studia dalších podkladů, odhadnout jeho odchylku oproti reálné vytíženosti.

1.1 Autoři a významné zdroje

Autory práce jsou Jakub Fukala a Adam Kozubek, kteří ji vytvořili jako projekt do předmětu modelování a simulace na FIT VUT v Brně. Významnými zdroji práce jsou zejména **výukové zdroje předmětu IMS** [8] [9].

Jako reprezentativní dataset vytížení byl pro naši práci použit reálný záznam denního vytížení z **Internet traffic archivu** [4]. Ten jsme dále zpracovali pomocí python knihoven.

Pro vytvoření relevantních modelů [11] autoscaleru, loadbalanceru a správné implementace obou přístupů škálování, jsme studovali internetové zdroje a knižní zdroj **The Datacenter as a Computer** [22]. Všechny použité zdroje jsou uvedeny na konci této dokumentace v sekci Bibliografie.

1.2 Prostředí a validace modelu

Model je vytvořen v jazyce **C++** a používá knihovnu **SIMLIB/C++** [8]. Ačkoliv je samotná simulace nedeterministická, její běh dá na různých strojích se stejnými vstupy srovnatelné výsledky, resp. nalezená odchylka je totožná s odchylkou [12] mezi více výstupy téže zopakované simulace na témže stroji.

Ačkoliv je model zjednodušený, má 15 nastavitelných parametrů a odpovídá realitě, protože byl vytvořen a nastaven na základě důkladné analýzy dostupných zdrojů. Také výstup simulace za daných podmínek byl srovnán s reálnými poznatky. Na základě toho model považujeme za validní [13].

2. Rozbor tématu a použitých metod a technologií

Moderní aplikace jsou stále častěji budovány na principu **mikroservisní architektury** [19], která umožňuje rozdělení aplikace na nezávislé komponenty – **mikroslužby**. Tato architektura umožňuje nejen nezávislý vývoj a nasazování jednotlivých částí systému, ale také jejich nezávislé škálování, čímž je dosaženo efektivního využití zdrojů. Klíčovou výzvou v prostředí mikroslužeb je **efektivní škálování**, které zajišťuje splnění dohodnuté úrovně služeb (**SLA**) [5] a zároveň minimalizaci nákladů na provoz infrastruktury.

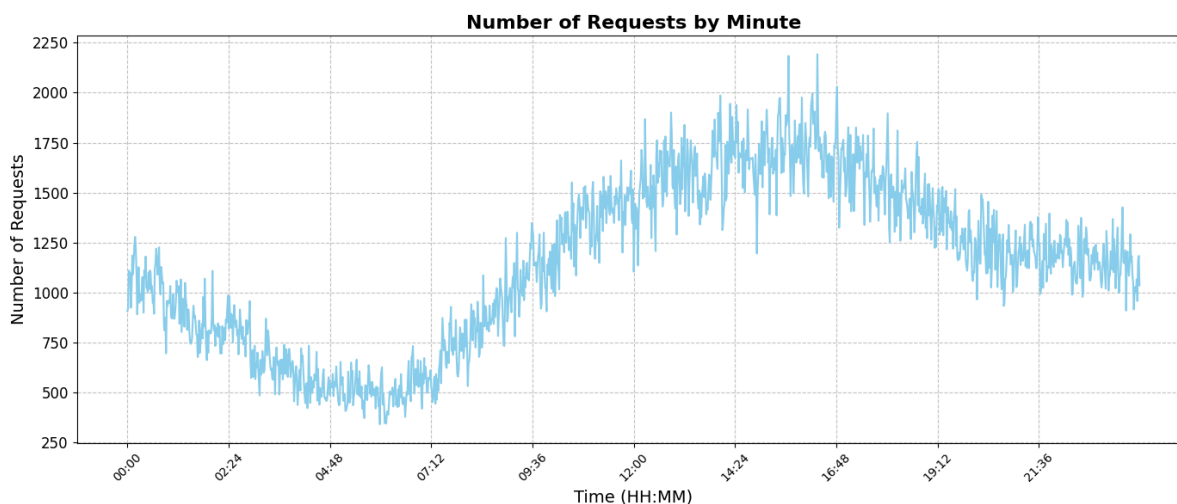
Pro náš model jsme hledali relevantní hodnoty parametrů pro klasickou webovou službu, na základě studií uvedených zdrojů jsme zvolili následující nastavení modelu:

- Typické webové aplikace mohou mít průměrné zatížení mezi **20 % až 60 % CPU** během běžného provozu, s krátkodobými špičkami až 90 % CPU při neočekávaných událostech [18].
- Vytvoření a zpřístupnění nového kontejneru na běžné platformě, jako je Kubernetes, trvá mezi **5 až 15 sekundami** v ideálním případě. Při zatížení systému se tento čas může zvýšit až na 30 sekund [18].
- Latence přesměrování: Dynamické algoritmy LoadBalanceru, jako Weighted Round Robin, mají přidanou latenci přesměrování obvykle **<10 ms** na požadavek. Tato hodnota se může zvýšit při přetížení na 30-50 ms [6].
- Směrodatná odchylka u odhadovaného zatížení (např. pro metriky CPU) se typicky pohybuje mezi **5 % až 25 %**, v závislosti na kvalitě dat a složitosti modelu [6].
- Reaktivní autoscalery reagují obvykle v časech **1-5 minut**, což zahrnuje čas na detekci potřeby škálování a spuštění nových instancí. [1]
- SLA jsme zvolili jako **95%** požadavků zpracovaných do **200 ms** [5].

Všechny tyto parametry společně s dalšími (podmínkami pro přeškálování, ...) jsou nastavitelné v sekci –PARAMETRY SIMULACE– tak, aby šla celá simulace spustit s modelem relevantním jiné konkrétní službě.

Zdrojem dat zátěže byl **Internet traffic archiv** [4], který nám poskytl reálná historická data počtu požadavků naměřených na konkrétním serveru. Naměřený vzorek dat obsahoval záznamy za sedm dní. Jednotlivé požadavky byly pro potřeby simulace shlukovány po jedné minutě a následně byla vytvořena průměrná hodnota požadavku v dané minutě. Výsledek byl vynesena do grafu o časové délce jednoho dne.

Zde je grafické znázornění vytvořené datové sady:



(Obrázek: 5 - Historická data zátěže jsou nahrána do prediktivního modelu.)

2.1. Popis použitých postupů pro vytvoření modelu a zdůvodnění

Simulace využívá modelování **škálování aplikací v kontejnerizačním prostředí** [21], kde je využito jak reaktivního, tak prediktivního přístupu k řízení počtu instancí. Reaktivní škálování reaguje na aktuální celkovou hodnotu zatížení kontejneru, zatímco prediktivní škálování využívá **historická data** a metody k předpovědi budoucí zátěže.

Postup modelování:

- **Generování požadavků:** Je dáno exponenciálním rozdělením [15] [2] [9]. Byly použity reálné datové soubory pro distribuci požadavků během 24 hodin [4].
- **Horizontální škálování:** Model škáluje počet instancí (kontejnerů) na základě průměrné zátěže (tj. průměrný počet současně odbavovaných požadavků). Doba spuštění nového kontejneru odpovídá hodnotám měřeným v prostředích Docker a Kubernetes [2].
- **Latence:** V modelu je zohledněno, že latence při obsluze požadavků narůstá při zvyšování zátěže daného kontejneru [22]. Pro zjednodušení jsme použili nejjednodušší model lineárního nárůstu [20].
- **Porovnání přístupů:** Simulace zahrnuje srovnání reaktivního a prediktivního přístupu. Prediktivní škálování dokáže díky předpovědím snížit riziko přetížení při náhlých nárůstu zátěže, zatímco reaktivní přístup může být efektivnější při minimalizaci nákladů ve stabilním prostředí [22].

Zdůvodnění postupů: Předpokládáme, že reaktivní škálování je jednodušší na implementaci a dobře vyhovuje scénářům s předvídatelnou a stabilní zátěží. Naopak prediktivní škálování vyžaduje více výpočetních zdrojů a kvalitní historická data [21], ale dokáže předcházet přetížení ve špičkách. Oba přístupy jsou v modelu implementovány a simulovány, aby byly vyhodnoceny jejich výhody a nevýhody v různých scénářích.

2.2. Popis původu použitých metod/technologií

Použité metody a technologie byly vybrány na základě ověřených vědeckých zdrojů a zkušeností z provozních systémů:

- **Kontejnerizace:** Simulace je navržena s ohledem na vlastnosti moderních kontejnerizačních nástrojů, jako je Docker a Kubernetes [2].
- **Škálování:** Pravidla škálování (např. prahové hodnoty) byla navržena na základě běžných praktik [1] a vlastních výsledků simulace.
- **Generování zátěže:** Datová sada obsahující záznamy o počtu požadavků za minutu pochází z reálných provozních měření a byla přizpůsobena pro simulaci pomocí normálního rozptylu s 15% odchylkou [4] [6].
- **LoadBlancing:** Jsme implementovali jednoduše tak, že nový požadavek je přiřazen kontejneru s nejnižším zatížením.

Tento přístup zajišťuje validitu modelu [13] a poskytuje užitečné srovnání škálovacích metod, které mohou být aplikovány v reálných systémech. Na základě těchto simulací lze posoudit efektivitu různých strategií škálování a jejich dopad na náklady i kvalitu služeb.

Pokročilé matematické funkce a funkce obsluhující běh simulace byly použity z knihovny **SIMLIB/C++** [15] [8], která je přesně na tento účel vytvořená.

3. Koncepce - modelářská témata

Náš model zanedbává dva aspekty. Prvním je náběh systému, kdy je připravený defaultně jediný kontejner a škálovací modely musí rychle zvýšit jejich počet, neboť zátěž je normální. Druhým je fakt, že latence od určitého bodu zátěže nenarůstá lineárně, respektive tento koncept nelze uvažovat, když chceme modelovat reálnou latenci.

Důvodem tohoto zjednodušení je fakt, že validita modelu není ovlivněna. Model ale není určen pro simulaci extrémních stavů, či náběhu a ukončování systému. Při náběhu systému se rychle dosáhne rovnováhy a v celkových statistikách se tento efekt neprojeví. Totéž platí o ukončování systému.

Dále se model omezuje jen na homogenní kontejnery, to znamená, že všechny kontejnery jsou považovány za identické z hlediska výkonu a latence, což zjednodušuje výpočty, ale odpovídá praxi v cloudovém prostředí s předdefinovanými instancemi.

Také jsou ignorovány síťové latence. Síťová latence je považována za zanedbatelnou, neboť většina aplikací má nízkou komunikaci mezi kontejnery v jednom uzlu.

Pro naše účely vyvození závěrů při běžném provozu systému běžné webové služby tak je model validní a dané závěry relevantní.

3.1 Způsob vyjádření konceptuálního modelu

Konceptuální model simuluje dynamické škálování aplikací v prostředí mikroservisní architektury. Tento model abstrahuje klíčové vlastnosti systému, jako je zátěž, odezva aplikací, přidávání a odebrání kontejnerů a jejich čas spuštění, a redukuje je na měřitelné a simulovatelné metriky. Model je vyjádřen kombinací schématického znázornění, stavového diagramu, Petriho sítí a matematických rovnic:

- **Schéma systému:** Na *obrázku 1* je znázorněn obecný tok požadavků od uživatele až po zpracování v rámci kontejnerů. Schéma ilustruje, jak zátěž vstupuje do systému, je rozdělena mezi aktivní kontejnery a třemi tečkami je naznačeno, jak probíhá škálování.

- **Matematické rovnice:**

- Lineární nárůst latence zpracování požadavku se zátěží kontejneru

$$doba_zpracování = min_doba_obsluhy \cdot (1 + \alpha \cdot zátěž_kontejneru)$$

- Odchylka reálných dat zátěže oproti predikci

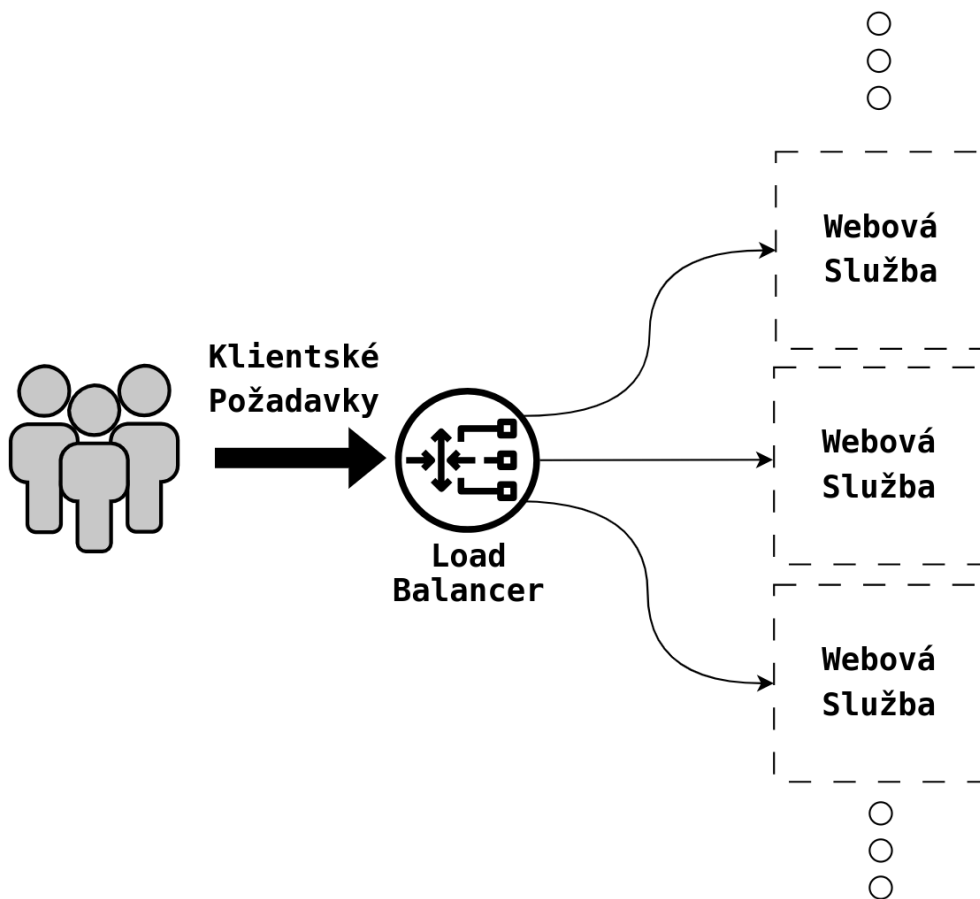
$$real_{LOAD} = Norm(predikce_{LOAD}; \sigma \cdot predikce_{LOAD})$$

- **Stavový diagram:** Na *obrázku 2* je znázorněn stavový diagram pro jednotlivé kontejnery, které mohou být v jedné z následujících stavů: **Neaktivní**, **Spouštění**, **Připravený** nebo **Aktivní**. Diagram také zahrnuje přechody mezi stavy, například

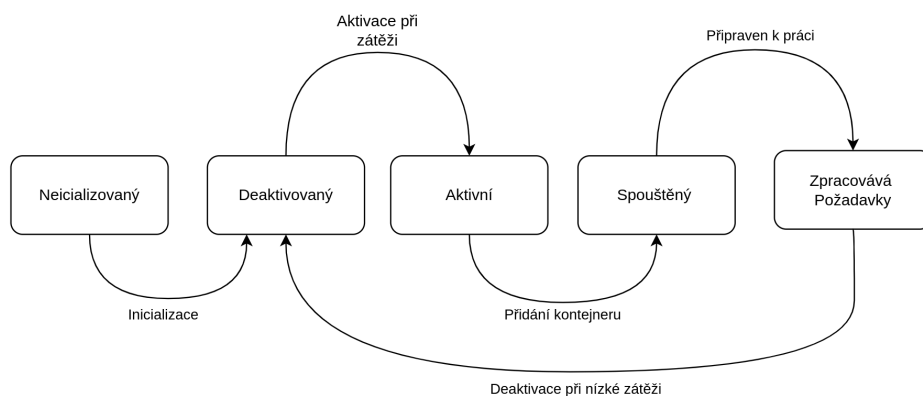
aktivaci při překročení zátěže nebo deaktivaci při nízkém vytížení.

- **Grafy Petriho Sítě [16]:** Na *obrázcích 3 a 4* jsou znázorněny Grafy Petriho Sítě [16]. Konkrétně na *obrázku 3* je část s LoadBalancerem a na *obrázku 4* je část s Autoscalerem. Tyto dvě části jsou propojené skrze místo Aktivní kontejnery, které je spravované AutoScalerem. Požadavky jsou generované s exponenciálním rozdělením, jehož středem je funkce $f(t)$, která odpovídá datasetu zatížení. Modely pro škálování zde jsou uvedeny pouze abstraktně v Autoscaleru v podmínkách přechodů pro přidání/odebrání kontejneru.

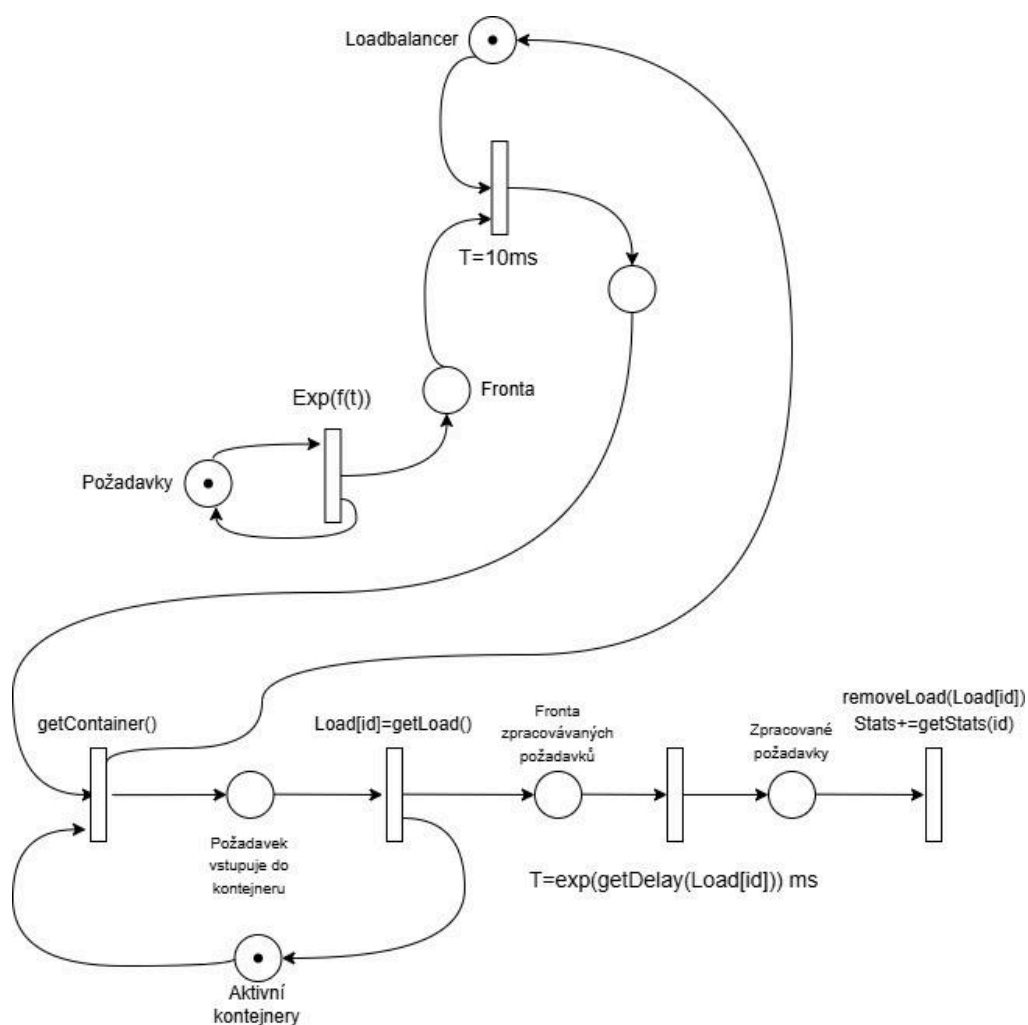
3.2 Forma vyjádření konceptuálního modelu



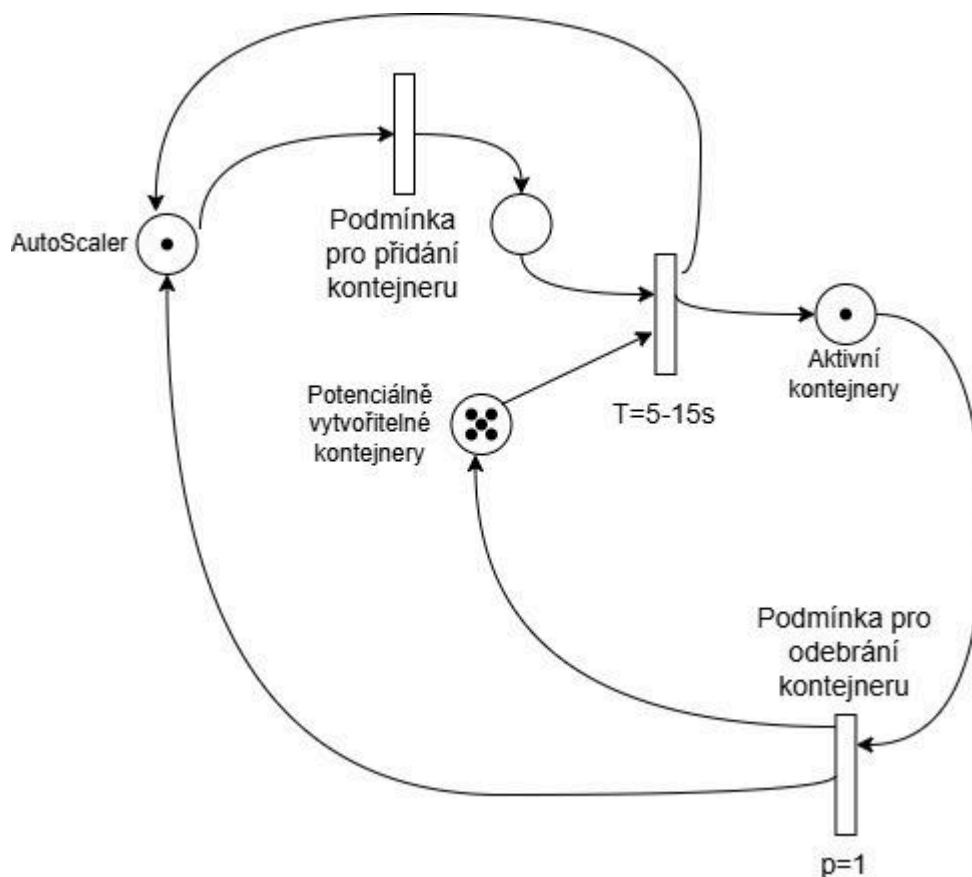
(Obrázek: 1 - Schéma systému)



(Obrázek: 2 - Stavový diagram kontejneru)



(Obrázek: 3 - Graf Petriho sítě [16] - LoadBalancer)



(Obrázek: 4 - Graf Petriho sítě [16] - AutoScaler)

3. Koncepce - implementační témata

Části kódu jsou blíže vysvětleny v kapitole 4. Konceptuální model je vysvětlen v předchozí kapitole. V této kapitole tedy jen popíšu základní informace o tom, jak funguje samotná implementace.

Na základě **datasetu zátěže** se s určitou kadencí pořád generují nově příchozí **požadavky**, které potřebují obsloužit. Na straně serveru jsou **kontejnery**, které je dokáží v určitém čase obsloužit. Tato doba je lineárně závislá na tom, jak je kontejner přetížený, podle vzorce: $t_{response} = t_{service} \cdot (1 + \alpha \cdot Load)$. Požadavek si při příchodu přiřadí nejméně zatížený kontejner (tj. LoadBalancing). Ten jej za $t_{response}$ zpracuje a uloží se statistiky. Zejména je kontrolováno, jestli $t_{response} > SLA_{time}$ kde **SLA** v defaultním nastavení kontroluje, že server zvládne 95% požadavků obsloužit do 200 ms.

Separátně jsou implementované **škálovací metody** (reaktivní a prediktivní). Škálování znamená, že mohou přidat nebo odebrat další kontejnery v definovaném konstantním čase. **Reaktivní přístup** se tak rozhoduje jen na základě aktuální vytíženosti,

určením maximálním a minimálním průměrným zatížením kontejnerů. Zatímco **prediktivní model** bere v potaz historický dataset a škálování řídí na základě požadované průměrné maximální procentuální zátěže kontejnerů.

Všechny tyto a další parametry, včetně výběru škálovací metody lze upravit a zvolit na začátku programu v sekci –PARAMETRY SIMULACE–.

Na konci simulace se vypíší statistiky, kromě SLA a průměrných vytížeností jednotlivých kontejnerů je důležitým údajem: **náklady na provoz**. Ty se vypočítají, jako součin ceny za provoz jednoho kontejneru a sumy celkového času provozu všech kontejnerů.

4. Architektura simulačního modelu/simulátoru

4.1 Mapování konceptuálního modelu do simulačního modelu

V této kapitole je podrobně popsáno, jak byl konceptuální model z kapitoly 3 implementován do konkrétního simulačního modelu za použití knihovny Simlib. Implementace využívá objektově orientované programování k modelování klíčových procesů a veličin. Níže jsou popsány jednotlivé prvky konceptuálního modelu a jejich odpovídající implementace v simulačním modelu.

Vytváření požadavků

- **Konceptuální model [17]:** Proces generuje příchozí požadavky podle dané zátěže modelované jako normální rozdělení na základě historických dat.
- **Implementace:**
 - **Třída RequestGenerator** dědí od `Event`. Aktivuje procesy požadavků v časových intervalech definovaných metodou `GetInterarrivalTime`, která využívá reálná data z pole `real_requests_per_minute`.
 - **Reálná data:** Data o zátěži jsou generována funkcí `GenerateRealRequestsPerMinute` a zohledňují variabilitu s odchylkou 15 %.

Požadavek

- **Konceptuální model [17]:** Požadavek čeká na přidělení volného kontejneru a je následně zpracován. Doba zpracování závisí na aktuální zátěži přiděleného kontejneru.

- **Implementace:**
 - **Třída Request** dědí od Process. Obsahuje metody pro simulaci přiřazení požadavku a zpracování v kontejneru.
 - **Výpočet doby zpracování:** Doba zpracování je modelována vztahem:

$$t_{response} = t_{service} \cdot (1 + \alpha \cdot Load)$$

kde α je koeficient ovlivňující nárůst latence při vyšší zátěži.

Kontejner

- **Konceptuální model [17]:** Kontejner má stavy **neaktivní**, **spouštějící**, **připravený** a **zpracovávající**. Sleduje zátěž, celkovou dobu aktivity a stav připravenosti.
- **Implementace:**
 - **Třída Container** reprezentuje jednotlivé kontejnery. Obsahuje atributy jako `is_ready`, `load` a `total_active_time`.
 - **Metody:**
 - **Activate a Deactivate:** Mění stav kontejneru.
 - **AcceptRequest a ReleaseRequest:** Přijímání a uvolnění požadavků.
 - **Start:** Nastavuje stav připravenosti po dokončení spouštění.
 - **Spouštění kontejneru:** Proces spouštění je reprezentován třídou `ContainerStartup`, která dědí od `Event`.

Škálování

- **Konceptuální model [17]:** Systém obsahuje dva různé škálovací algoritmy:
 - **Reaktivní škálování** reaguje na aktuální zátěž.
 - **Prediktivní škálování** využívá historická data k predikci budoucí zátěže.
- **Implementace:**
 - **Reaktivní metoda:**
 - **Třída ReactiveAutoscaler:** V pravidelných intervalech kontroluje průměrnou zátěž aktivních kontejnerů a přidává nebo odstraňuje kontejnery podle prahových hodnot (`SCALE_UP_LOAD` a `SCALE_DOWN_LOAD`).
 - **Prediktivní metoda:**
 - **Třída PredictiveAutoscaler:** Predikuje maximální počet požadavků za interval a vypočítává potřebný počet kontejnerů na základě požadované zátěže a latence.

Metody škálování

- Výpočet počtu kontejnerů:

- **Reaktivní škálování:**

$$Load_{average} = \frac{Total\ Load}{Active\ Containers}$$

Pokud $Load_{average}$ překročí $SCALE_UP_LOAD$, přidají se kontejnery, naopak při poklesu pod $SCALE_DOWN_LOAD$ se odstraní.

- **Prediktivní škálování:**

$$N_{containers} = \frac{R_{predicted} \cdot t_{response}}{Load_{desired}}$$

Predikuje se vždy dopředu o dobu, do níž nebude možno provést další predikci a spustit/odebrat další kontejnery.

Statistiky

- **Konceptuální model:** Statistické výstupy zahrnují dobu odezvy, porušení SLA a náklady na provoz.
- **Implementace:**
 - **Knihovna Simlib** poskytuje třídy `Stat` a `Histogram`, které ukládají data a generují výstupy.
 - Statistika doby odezvy je implementována jako instance `Stat response_time_stat` a `Histogram response_time_hist`.

4.2 Mapování abstraktního modelu na implementaci

Konceptuální prvek	Třída/metoda v implementaci
Požadavky	<code>Request</code> , <code>RequestGenerator</code>
Kontejnery	<code>Request</code> , <code>RequestGenerator</code>
Reaktivní škálování	<code>ReactiveAutoscaler</code>
Prediktivní škálování	<code>PredictiveAutoscaler</code>
Generování reálné zátěže	<code>GenerateRealRequestsPerMinute</code>
Statistika doby odezvy	<code>GeneratePredictedLoad</code>
Náklady na provoz	Výpočet v hlavní funkci <code>main</code>

5. Podstata simulačních experimentů a jejich průběh

5.1 Postup experimentování

Stejně nastavení simulace jsme vždy vyzkoušeli pro prediktivní i reaktivní model, poté jsme formulovali pozorování. Postupně jsme experimentovali s různými nastavením parametrů a na základě porovnání výsledků simulačních běhů, formulovali závěry.

5.2 Dokumentace jednotlivých experimentů

Výsledek simulace pro defaultní nastavení, relevantní modelové webové aplikace:

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102

Max = 1.604

Number of records = 41247850

Average value = 0.144613

Standard deviation = 0.0354822

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	29511807	0.715475	0.715475
0.150	0.200	11290560	0.273725	0.989200
0.200	0.250	325089	0.007881	0.997081
0.250	0.300	53699	0.001302	0.998383
0.300	0.350	16145	0.000391	0.998774
0.350	0.400	13778	0.000334	0.999109
0.400	0.450	4577	0.000111	0.999219
0.450	0.500	1343	0.000033	0.999252
0.500	0.550	965	0.000023	0.999275
0.550	0.600	293	0.000007	0.999283
0.600	0.650	295	0.000007	0.999290
0.650	0.700	231	0.000006	0.999295
0.700	0.750	536	0.000013	0.999308
0.750	0.800	407	0.000010	0.999318
0.800	0.850	1025	0.000025	0.999343
0.850	0.900	926	0.000022	0.999365
0.900	0.950	739	0.000018	0.999383
0.950	1.000	739	0.000018	0.999401

SLA splněno pro 98.92% požadavků.

Kontejner 0 průměrná zátěž: 24.4363

Kontejner 1 průměrná zátěž: 21.1257

Kontejner 2 průměrná zátěž: 20.9315

Kontejner 3 průměrná zátěž: 20.1985

Kontejner 4 průměrná zátěž: 19.0198

Celkové náklady na provoz: 8.28542

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102

Max = 0.808

Number of records = 34208612

Average value = 0.164025

Standard deviation = 0.0254202

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	9966483	0.291344	0.291344
0.150	0.200	21800608	0.637284	0.928628
0.200	0.250	2135220	0.062418	0.991046
0.250	0.300	280443	0.008198	0.999244
0.300	0.350	19609	0.000573	0.999817
0.350	0.400	180	0.000005	0.999823
0.400	0.450	293	0.000009	0.999831
0.450	0.500	452	0.000013	0.999844
0.500	0.550	328	0.000010	0.999854
0.550	0.600	343	0.000010	0.999864
0.600	0.650	211	0.000006	0.999870
0.650	0.700	566	0.000017	0.999887
0.700	0.750	1343	0.000039	0.999926
0.750	0.800	2503	0.000073	0.999999
0.800	0.850	30	0.000001	1.000000
0.850	0.900	0	0.000000	1.000000
0.900	0.950	0	0.000000	1.000000
0.950	1.000	0	0.000000	1.000000

SLA splněno pro 92.8628% požadavků.

Kontejner 0 průměrná zátěž: 32.0044

Kontejner 1 průměrná zátěž: 30.776

Kontejner 2 průměrná zátěž: 32.7931

Kontejner 3 průměrná zátěž: 28.7082

Celkové náklady na provoz: 6.51583

Pozorování:

Prediktivní přístup ušetřil 20% nákladů a vystačil si jen se 4 kontejnery oproti 5 u reaktivního přístupu, zato ale nesplnil SLA. Zjišťujeme, že požadované 85% vytížení kontejnerů u prediktivního modelu je příliš vysoké.

Nastavení: Požadované vytížení kontejneru: DESIRED_PERCENTAGE_LOAD = 75%

Reaktivní Škálování

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102Max = 1.604

Number of records = 41247850

Average value = 0.144613

Standard deviation = 0.0354822

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	29511807	0.715475	0.715475
0.150	0.200	11290560	0.273725	0.989200
0.200	0.250	325089	0.007881	0.997081
0.250	0.300	53699	0.001302	0.998383
0.300	0.350	16145	0.000391	0.998774
0.350	0.400	13778	0.000334	0.999109
0.400	0.450	4577	0.000111	0.999219
0.450	0.500	1343	0.000033	0.999252
0.500	0.550	965	0.000023	0.999275
0.550	0.600	293	0.000007	0.999283
0.600	0.650	295	0.000007	0.999290
0.650	0.700	231	0.000006	0.999295
0.700	0.750	536	0.000013	0.999308
0.750	0.800	407	0.000010	0.999318
0.800	0.850	1025	0.000025	0.999343
0.850	0.900	926	0.000022	0.999365
0.900	0.950	739	0.000018	0.999383
0.950	1.000	739	0.000018	0.999401

SLA splněno pro 98.92% požadavků.
Kontejner 0 průměrná zátěž: 24.4363
Kontejner 1 průměrná zátěž: 21.1257
Kontejner 2 průměrná zátěž: 20.9315
Kontejner 3 průměrná zátěž: 20.1985
Kontejner 4 průměrná zátěž: 19.0198
Celkové náklady na provoz: 8.28542

Prediktivní Škálování

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102Max = 1.19

Number of records = 41379611

Average value = 0.158067

Standard deviation = 0.022999

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	15500216	0.374586	0.374586
0.150	0.200	24342823	0.588281	0.962866
0.200	0.250	1462189	0.035336	0.998202
0.250	0.300	28952	0.000700	0.998902
0.300	0.350	34456	0.000833	0.999735
0.350	0.400	4328	0.000105	0.999839
0.400	0.450	130	0.000003	0.999843
0.450	0.500	111	0.000003	0.999845
0.500	0.550	117	0.000003	0.999848
0.550	0.600	272	0.000007	0.999855
0.600	0.650	239	0.000006	0.999860
0.650	0.700	200	0.000005	0.999865
0.700	0.750	441	0.000011	0.999876
0.750	0.800	489	0.000012	0.999888
0.800	0.850	188	0.000005	0.999892
0.850	0.900	187	0.000005	0.999897
0.900	0.950	806	0.000019	0.999916
0.950	1.000	1026	0.000025	0.999941

SLA splněno pro 96.2866% požadavků.
Kontejner 0 průměrná zátěž: 28.91
Kontejner 1 průměrná zátěž: 28.0601
Kontejner 2 průměrná zátěž: 28.8112
Kontejner 3 průměrná zátěž: 28.0335
Kontejner 4 průměrná zátěž: 30.6958
Celkové náklady na provoz: 6.85333

Pozorování:

Prediktivní přístup vesměs dokázal udržet maximální průměrné vytížení všech kontejnerů přibližně na 75%, tak že splnil SLA. Už musel aktivovat i 5. kontejner, ale všechny kontejnery využíval efektivněji než reaktivní přístup. Tudíž prediktivní přístup stále ušetřil 17% nákladů.

Abychom se pokusili s reaktivním přístupem dohnat nízkou cenu provozu prediktivního škálování, musíme snížit spotřebu kontejnerů.

Nastavení: Nižší spotřeba kontejnerů reaktivního přístupu.

SCALE_DOWN_LOAD_PERCENTAGE = 40

Reaktivní Škálování

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102Max = 2.82
Number of records = 41697443
Average value = 0.169814
Standard deviation = 0.0767422

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	15332576	0.367710	0.367710
0.150	0.200	21016045	0.504013	0.871723
0.200	0.250	3747109	0.089864	0.961587
0.250	0.300	969290	0.023246	0.984833
0.300	0.350	269081	0.006453	0.991286
0.350	0.400	92511	0.002219	0.993505
0.400	0.450	40849	0.000980	0.994485
0.450	0.500	29061	0.000697	0.995181
0.500	0.550	33436	0.000802	0.995983
0.550	0.600	24454	0.000586	0.996570
0.600	0.650	18421	0.000442	0.997012
0.650	0.700	16689	0.000400	0.997412
0.700	0.750	14891	0.000357	0.997769
0.750	0.800	13750	0.000330	0.998099
0.800	0.850	5354	0.000128	0.998227
0.850	0.900	5672	0.000136	0.998363
0.900	0.950	5119	0.000123	0.998486
0.950	1.000	1916	0.000046	0.998532

SLA splněno pro 87.1723% požadavků.
Kontejner 0 průměrná zátěž: 42.0124
Kontejner 1 průměrná zátěž: 32.1813
Kontejner 2 průměrná zátěž: 28.6484
Kontejner 3 průměrná zátěž: 25.3805
Kontejner 4 průměrná zátěž: 20.89
Celkové náklady na provoz: 6.67625

Prediktivní Škálování

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102Max = 1.19
Number of records = 41379611
Average value = 0.158067
Standard deviation = 0.022999

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	15500216	0.374586	0.374586
0.150	0.200	24342823	0.588281	0.962866
0.200	0.250	1462189	0.035336	0.998202
0.250	0.300	28952	0.000700	0.998902
0.300	0.350	34456	0.000833	0.999735
0.350	0.400	4328	0.000105	0.999839
0.400	0.450	130	0.000003	0.999843
0.450	0.500	111	0.000003	0.999845
0.500	0.550	117	0.000003	0.999848
0.550	0.600	272	0.000007	0.999855
0.600	0.650	239	0.000006	0.999860
0.650	0.700	200	0.000005	0.999865
0.700	0.750	441	0.000011	0.999876
0.750	0.800	489	0.000012	0.999888
0.800	0.850	188	0.000005	0.999892
0.850	0.900	187	0.000005	0.999897
0.900	0.950	806	0.000019	0.999916
0.950	1.000	1026	0.000025	0.999941

SLA splněno pro 96.2866% požadavků.
Kontejner 0 průměrná zátěž: 28.91
Kontejner 1 průměrná zátěž: 28.0601
Kontejner 2 průměrná zátěž: 28.8112
Kontejner 3 průměrná zátěž: 28.0335
Kontejner 4 průměrná zátěž: 30.6958
Celkové náklady na provoz: 6.85333

Pozorování: Abychom se dostali na onu nižší cenu, museli jsme zvýšit spodní hranici pro přeskálování dolů na 40%. Reaktivní model ale při takovém nastavení nezvládl SLA.

Nastavení: Větší rozptýlení zátěže. STANDARD_DEVIATION = 0.25

Reaktivní Škálování

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102Max = 1.834

Number of records = 41414014

Average value = 0.14929

Standard deviation = 0.0455508

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	26205743	0.632775	0.632775
0.150	0.200	13593388	0.328232	0.961006
0.200	0.250	1222211	0.029512	0.990518
0.250	0.300	273904	0.006614	0.997132
0.300	0.350	35541	0.000858	0.997990
0.350	0.400	7057	0.000170	0.998161
0.400	0.450	14643	0.000354	0.998514
0.450	0.500	7326	0.000177	0.998691
0.500	0.550	7092	0.000171	0.998862
0.550	0.600	7187	0.000174	0.999036
0.600	0.650	4862	0.000117	0.999153
0.650	0.700	1058	0.000026	0.999179
0.700	0.750	199	0.000005	0.999184
0.750	0.800	248	0.000006	0.999190
0.800	0.850	394	0.000010	0.999199
0.850	0.900	1303	0.000031	0.999231
0.900	0.950	1185	0.000029	0.999259
0.950	1.000	695	0.000017	0.999276

SLA splněno pro 96.1007% požadavků.

Kontejner 0 průměrná zátěž: 28.2624

Kontejner 1 průměrná zátěž: 23.3422

Kontejner 2 průměrná zátěž: 22.6131

Kontejner 3 průměrná zátěž: 21.4859

Kontejner 4 průměrná zátěž: 20.2378

Kontejner 5 průměrná zátěž: 18.693

Kontejner 6 průměrná zátěž: 12.0372

Celkové náklady na provoz: 8.29458

Prediktivní Škálování

HISTOGRAM Histogram doby odezvy

STATISTIC

Min = 0.102Max = 0.494

Number of records = 41090551

Average value = 0.162798

Standard deviation = 0.0306001

from	to	n	rel	sum
0.000	0.050	0	0.000000	0.000000
0.050	0.100	0	0.000000	0.000000
0.100	0.150	15278504	0.371825	0.371825
0.150	0.200	21808646	0.530746	0.902571
0.200	0.250	3477784	0.084637	0.987208
0.250	0.300	380215	0.009253	0.996461
0.300	0.350	60748	0.001478	0.997940
0.350	0.400	25778	0.000627	0.998567
0.400	0.450	53142	0.001293	0.999860
0.450	0.500	5734	0.000140	1.000000
0.500	0.550	0	0.000000	1.000000
0.550	0.600	0	0.000000	1.000000
0.600	0.650	0	0.000000	1.000000
0.650	0.700	0	0.000000	1.000000
0.700	0.750	0	0.000000	1.000000
0.750	0.800	0	0.000000	1.000000
0.800	0.850	0	0.000000	1.000000
0.850	0.900	0	0.000000	1.000000
0.900	0.950	0	0.000000	1.000000
0.950	1.000	0	0.000000	1.000000

SLA splněno pro 90.2571% požadavků.

Kontejner 0 průměrná zátěž: 31.2087

Kontejner 1 průměrná zátěž: 30.4776

Kontejner 2 průměrná zátěž: 31.3833

Kontejner 3 průměrná zátěž: 29.8671

Kontejner 4 průměrná zátěž: 36.118

Celkové náklady na provoz: 6.85333

Pozorování: Významné rozptýlení minutové zátěže způsobilo prudké změny v zátěži a tudíž prediktivní model nesplnil SLA a reaktivní jen těsně. Cenový rozdíl zůstává přibližně stejný. Zvýšení rozptýlení přineslo snížení přesnosti predikce a tudíž selhání.

Hypotéza: Pokud by si byl prediktivní model vědom své odchylky predikce od reality, snížil své požadované vytížení kontejneru, mohl by rozdíl vyrovnat.

Nastavení: DESIRED_PERCENTAGE_LOAD = 65

Reaktivní Škálování	Prediktivní Škálování																																																																																																																																																																																																																		
<div><div>HISTOGRAM Histogram doby odezvy</div><div><div>STATISTIC</div><div>Min = 0.102Max = 1.834 Number of records = 41414014 Average value = 0.14929 Standard deviation = 0.0455508</div></div><table><tr><th>from</th><th>to</th><th>n</th><th>rel</th><th>sum</th></tr><tr><td>0.000</td><td>0.050</td><td>0</td><td>0.000000</td><td>0.000000</td></tr><tr><td>0.050</td><td>0.100</td><td>0</td><td>0.000000</td><td>0.000000</td></tr><tr><td>0.100</td><td>0.150</td><td>26205743</td><td>0.632775</td><td>0.632775</td></tr><tr><td>0.150</td><td>0.200</td><td>13593388</td><td>0.328232</td><td>0.961006</td></tr><tr><td>0.200</td><td>0.250</td><td>1222211</td><td>0.029512</td><td>0.990518</td></tr><tr><td>0.250</td><td>0.300</td><td>273904</td><td>0.006614</td><td>0.997132</td></tr><tr><td>0.300</td><td>0.350</td><td>35541</td><td>0.000858</td><td>0.997990</td></tr><tr><td>0.350</td><td>0.400</td><td>7057</td><td>0.000170</td><td>0.998161</td></tr><tr><td>0.400</td><td>0.450</td><td>14643</td><td>0.000354</td><td>0.998514</td></tr><tr><td>0.450</td><td>0.500</td><td>7326</td><td>0.000177</td><td>0.998691</td></tr><tr><td>0.500</td><td>0.550</td><td>7092</td><td>0.000171</td><td>0.998862</td></tr><tr><td>0.550</td><td>0.600</td><td>7187</td><td>0.000174</td><td>0.999036</td></tr><tr><td>0.600</td><td>0.650</td><td>4862</td><td>0.000117</td><td>0.999153</td></tr><tr><td>0.650</td><td>0.700</td><td>1058</td><td>0.000026</td><td>0.999179</td></tr><tr><td>0.700</td><td>0.750</td><td>199</td><td>0.000005</td><td>0.999184</td></tr><tr><td>0.750</td><td>0.800</td><td>248</td><td>0.000006</td><td>0.999190</td></tr><tr><td>0.800</td><td>0.850</td><td>394</td><td>0.000010</td><td>0.999199</td></tr><tr><td>0.850</td><td>0.900</td><td>1303</td><td>0.000031</td><td>0.999231</td></tr><tr><td>0.900</td><td>0.950</td><td>1185</td><td>0.000029</td><td>0.999259</td></tr><tr><td>0.950</td><td>1.000</td><td>695</td><td>0.000017</td><td>0.999276</td></tr></table><div>SLA splněno pro 96.1007% požadavků. Kontejner 0 průměrná zátěž: 28.2624 Kontejner 1 průměrná zátěž: 23.3422 Kontejner 2 průměrná zátěž: 22.6131 Kontejner 3 průměrná zátěž: 21.4859 Kontejner 4 průměrná zátěž: 20.2378 Kontejner 5 průměrná zátěž: 18.693 Kontejner 6 průměrná zátěž: 12.0372 Celkové náklady na provoz: 8.29458</div></div>	from	to	n	rel	sum	0.000	0.050	0	0.000000	0.000000	0.050	0.100	0	0.000000	0.000000	0.100	0.150	26205743	0.632775	0.632775	0.150	0.200	13593388	0.328232	0.961006	0.200	0.250	1222211	0.029512	0.990518	0.250	0.300	273904	0.006614	0.997132	0.300	0.350	35541	0.000858	0.997990	0.350	0.400	7057	0.000170	0.998161	0.400	0.450	14643	0.000354	0.998514	0.450	0.500	7326	0.000177	0.998691	0.500	0.550	7092	0.000171	0.998862	0.550	0.600	7187	0.000174	0.999036	0.600	0.650	4862	0.000117	0.999153	0.650	0.700	1058	0.000026	0.999179	0.700	0.750	199	0.000005	0.999184	0.750	0.800	248	0.000006	0.999190	0.800	0.850	394	0.000010	0.999199	0.850	0.900	1303	0.000031	0.999231	0.900	0.950	1185	0.000029	0.999259	0.950	1.000	695	0.000017	0.999276	<div><div>HISTOGRAM Histogram doby odezvy</div><div><div>STATISTIC</div><div>Min = 0.102Max = 0.348 Number of records = 41239625 Average value = 0.155389 Standard deviation = 0.0240711</div></div><table><tr><th>from</th><th>to</th><th>n</th><th>rel</th><th>sum</th></tr><tr><td>0.000</td><td>0.050</td><td>0</td><td>0.000000</td><td>0.000000</td></tr><tr><td>0.050</td><td>0.100</td><td>0</td><td>0.000000</td><td>0.000000</td></tr><tr><td>0.100</td><td>0.150</td><td>19228235</td><td>0.466256</td><td>0.466256</td></tr><tr><td>0.150</td><td>0.200</td><td>20030243</td><td>0.485704</td><td>0.951960</td></tr><tr><td>0.200</td><td>0.250</td><td>1753710</td><td>0.042525</td><td>0.994485</td></tr><tr><td>0.250</td><td>0.300</td><td>223174</td><td>0.005412</td><td>0.999897</td></tr><tr><td>0.300</td><td>0.350</td><td>4263</td><td>0.000103</td><td>1.000000</td></tr><tr><td>0.350</td><td>0.400</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.400</td><td>0.450</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.450</td><td>0.500</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.500</td><td>0.550</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.550</td><td>0.600</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.600</td><td>0.650</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.650</td><td>0.700</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.700</td><td>0.750</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.750</td><td>0.800</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.800</td><td>0.850</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.850</td><td>0.900</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.900</td><td>0.950</td><td>0</td><td>0.000000</td><td>1.000000</td></tr><tr><td>0.950</td><td>1.000</td><td>0</td><td>0.000000</td><td>1.000000</td></tr></table><div>SLA splněno pro 95.196% požadavků. Kontejner 0 průměrná zátěž: 27.3119 Kontejner 1 průměrná zátěž: 26.8666 Kontejner 2 průměrná zátěž: 27.6254 Kontejner 3 průměrná zátěž: 27.1301 Kontejner 4 průměrná zátěž: 23.7846 Celkové náklady na provoz: 7.34583</div></div>	from	to	n	rel	sum	0.000	0.050	0	0.000000	0.000000	0.050	0.100	0	0.000000	0.000000	0.100	0.150	19228235	0.466256	0.466256	0.150	0.200	20030243	0.485704	0.951960	0.200	0.250	1753710	0.042525	0.994485	0.250	0.300	223174	0.005412	0.999897	0.300	0.350	4263	0.000103	1.000000	0.350	0.400	0	0.000000	1.000000	0.400	0.450	0	0.000000	1.000000	0.450	0.500	0	0.000000	1.000000	0.500	0.550	0	0.000000	1.000000	0.550	0.600	0	0.000000	1.000000	0.600	0.650	0	0.000000	1.000000	0.650	0.700	0	0.000000	1.000000	0.700	0.750	0	0.000000	1.000000	0.750	0.800	0	0.000000	1.000000	0.800	0.850	0	0.000000	1.000000	0.850	0.900	0	0.000000	1.000000	0.900	0.950	0	0.000000	1.000000	0.950	1.000	0	0.000000	1.000000
from	to	n	rel	sum																																																																																																																																																																																																															
0.000	0.050	0	0.000000	0.000000																																																																																																																																																																																																															
0.050	0.100	0	0.000000	0.000000																																																																																																																																																																																																															
0.100	0.150	26205743	0.632775	0.632775																																																																																																																																																																																																															
0.150	0.200	13593388	0.328232	0.961006																																																																																																																																																																																																															
0.200	0.250	1222211	0.029512	0.990518																																																																																																																																																																																																															
0.250	0.300	273904	0.006614	0.997132																																																																																																																																																																																																															
0.300	0.350	35541	0.000858	0.997990																																																																																																																																																																																																															
0.350	0.400	7057	0.000170	0.998161																																																																																																																																																																																																															
0.400	0.450	14643	0.000354	0.998514																																																																																																																																																																																																															
0.450	0.500	7326	0.000177	0.998691																																																																																																																																																																																																															
0.500	0.550	7092	0.000171	0.998862																																																																																																																																																																																																															
0.550	0.600	7187	0.000174	0.999036																																																																																																																																																																																																															
0.600	0.650	4862	0.000117	0.999153																																																																																																																																																																																																															
0.650	0.700	1058	0.000026	0.999179																																																																																																																																																																																																															
0.700	0.750	199	0.000005	0.999184																																																																																																																																																																																																															
0.750	0.800	248	0.000006	0.999190																																																																																																																																																																																																															
0.800	0.850	394	0.000010	0.999199																																																																																																																																																																																																															
0.850	0.900	1303	0.000031	0.999231																																																																																																																																																																																																															
0.900	0.950	1185	0.000029	0.999259																																																																																																																																																																																																															
0.950	1.000	695	0.000017	0.999276																																																																																																																																																																																																															
from	to	n	rel	sum																																																																																																																																																																																																															
0.000	0.050	0	0.000000	0.000000																																																																																																																																																																																																															
0.050	0.100	0	0.000000	0.000000																																																																																																																																																																																																															
0.100	0.150	19228235	0.466256	0.466256																																																																																																																																																																																																															
0.150	0.200	20030243	0.485704	0.951960																																																																																																																																																																																																															
0.200	0.250	1753710	0.042525	0.994485																																																																																																																																																																																																															
0.250	0.300	223174	0.005412	0.999897																																																																																																																																																																																																															
0.300	0.350	4263	0.000103	1.000000																																																																																																																																																																																																															
0.350	0.400	0	0.000000	1.000000																																																																																																																																																																																																															
0.400	0.450	0	0.000000	1.000000																																																																																																																																																																																																															
0.450	0.500	0	0.000000	1.000000																																																																																																																																																																																																															
0.500	0.550	0	0.000000	1.000000																																																																																																																																																																																																															
0.550	0.600	0	0.000000	1.000000																																																																																																																																																																																																															
0.600	0.650	0	0.000000	1.000000																																																																																																																																																																																																															
0.650	0.700	0	0.000000	1.000000																																																																																																																																																																																																															
0.700	0.750	0	0.000000	1.000000																																																																																																																																																																																																															
0.750	0.800	0	0.000000	1.000000																																																																																																																																																																																																															
0.800	0.850	0	0.000000	1.000000																																																																																																																																																																																																															
0.850	0.900	0	0.000000	1.000000																																																																																																																																																																																																															
0.900	0.950	0	0.000000	1.000000																																																																																																																																																																																																															
0.950	1.000	0	0.000000	1.000000																																																																																																																																																																																																															

Pozorování: Prediktivní model za těchto okolností též zvládnul SLA, ušetřil ale jen polovinu nákladů co před tím.

6. Shrnutí simulačních experimentů a závěr

V rámci tohoto projektu vznikl nástroj, který umožňuje testovat hypotézy, porovnat modely a simulovat provoz serveru při daných podmínkách.

Průběžné kontrolní výpisy simulace jsou smysluplné a pomohou porozumět, co se při simulaci odehrává.

Při zkoušení různých nastavení simulace odpovídal její výsledek předpokladům, nabytých na základě studia obecných závěrů použitých zdrojů. Validita modelu tak byla ověřena.

Oproti očekáváním dokázal prediktivní přístup snížit cenu za provoz i v relativně rovnoměrném provozu. Jako slabina tohoto přístupu se ukázala přesnost predikce. Když klesla, a směrodatná odchylka oproti reálu dosahovala 25%, přístup nebyl schopen splnit SLA, popřípadě rychle zvyšoval náklady při snaze toho dosáhnout.

Naproti tomu reaktivní model dokázal při bezpečném nastavení (20% - 70% průměrné zátěže) splnit SLA vždy. Nevýhodou tohoto přístupu je, že nedokáže v předvídatelné zátěži cenově konkurovat prediktivnímu modelu ani při progresivnějším nastavení.

7. Bibliografie

- [1] ALHARTHI, Saleha; ALSHAMSI, Afra; ALSEIARI, Anoud; ALWARAFY Abdulmalik; 2024. Auto-Scaling Techniques in Cloud Computing: Issues and Research Directions. Dostupné z: <https://www.mdpi.com/1424-8220/24/17/5551>. [cit. 2024-11-30].
- [2] CLOUD NATIVE COMPUTING FOUNDATION; 2024. Pod Lifecycle. Dostupné z: <https://kubernetes.io/docs/concepts/workloads/pods/pod-lifecycle/>. [cit. 2024-11-30].
- [3] CONRAN, Matt; 2015. Load Balancing and Scale-Out Architectures. Dostupné z: <https://network-insight.net/2015/02/26/load-balancing-and-scale-out-architectures/>. [cit. 2024-11-29].
- [4] DANZIG, Peter; et al. 2008. The Internet Traffic Archive. Dostupné z: <https://ita.ee.lbl.gov/html/traces.html>. [cit. 2024-11-29].
- [5] Chris Jones, John Wilkes, and Niall Murphy with Cody Smith; 2017. Service Level Objectives. Dostupné z: <https://sre.google/sre-book/service-level-objectives/>. [cit. 2024-11-30].
- [6] NEWMAN, David; 2024. Load Balancing Fundamentals: How Load Balancers Work. Dostupné z: <https://www.linode.com/docs/guides/load-balancing-fundamentals/>. [cit. 2024-11-27].
- [7] PERINGER, Petr; MARTINEK, David; LEŠKA, David; 2021. SIMLIB/C++ Documentation. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/SIMLIB/doc/html/>. [cit. 2024-11-26].
- [8] PERINGER, Petr; 2022. Modelování a simulace - Studijní opora. Dostupné z: https://moodle.vut.cz/pluginfile.php/900581/mod_resource/content/2/opora-ims.pdf. [cit. 2024-11-28].
- [9] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf>. [cit. 2024-11-26].
- [10] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 33. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [11] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 44. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [12] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 82. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [13] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 37. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [14] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 91. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [15] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 38. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [16] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 127. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].
- [17] PERINGER, Petr; HRUBÝ Martin; 2024. Modelování a simulace - Prezentace slice č. 48. Dostupné z: <https://www.fit.vut.cz/person/peringer/public/IMS/prednasky/IMS.pdf> [cit. 2024-11-26].

- [18] RABIU, Shamsuddeen; CHAN, Huah Yong; 2022. A Cloud-Based Container Microservices: A Review on Load-Balancing and Auto-Scaling Issues. Dostupné z: https://www.researchgate.net/publication/366562529_A_Cloud-Based_Container_Microservices_A_Review_on_Load-Balancing_and_Auto-Scaling_Issues. [cit. 2024-11-27].
- [19] MICROSOFT - AzureDevOps; 2000. Microservices architecture design. Dostupné z: <https://learn.microsoft.com/en-us/azure/architecture/microservices/>. [cit. 2024-11-29].
- [20] NIST; Technical Series Publications. Dostupné z: <https://nvlpubs.nist.gov/>. [cit. 2024-11-30].
- [21] SAFA, Bouguezz; 2024. Název: Podnázev. Dostupné z: <https://www.baeldung.com/cs/scaling-horizontally-vertically>. [cit. 2024-11-30].
- [22] BARROSO, Luiz André; HÖLZLE, Urs; RANGANATHAN, Parthasarathy; 2019. The datacenter as a computer. Dostupné z: <https://link.springer.com/book/10.1007/978-3-031-01761-2>. [cit. 2024-11-29].