

ارائه رویکردی نوین مبتنی بر یادگیری بیزی برای تامین کارآمد منابع برنامه‌های کاربردی در محیط‌های ابری

مصطفی قبائی آرانی*

گروه مهندسی کامپیوتر، واحد قم، دانشگاه آزاد اسلامی، قم، ایران
پست الکترونیکی: m.ghobaei@qom-iau.ac.ir

سمانه کربلایی مهدی

گروه مهندسی کامپیوتر، واحد محلات، دانشگاه آزاد اسلامی، مرکزی، ایران
پست الکترونیکی: samaneh.karbalaie@gmail.com

چکیده

میانگین بهره‌وری، کاهش ۳ درصدی زمان پاسخ‌دهی، در نتیجه کاهش ۵ درصدی هزینه تمام شده و افزایش ۱ درصدی سود حاصل شده است. واژه‌های کلیدی: رایانش ابری، برنامه‌های چندلایه، مقیاس‌پذیری، رگرسیون خطی^۲، نظریه بیز^۳.

مقدمه

گسترش اینترنت و فراگیری استفاده از آن، تاثیر شگرفی بر زندگی بشر داشته و خدماتی که از طریق اینترنت می‌توان ارائه داد دائماً در حال تغییر است. با پیشرفت فناوری اطلاعات نیاز به انجام کارهای محاسباتی در همه جا و همه مکان به وجود آمده است. همچنین نیاز به این هست که افراد بتوانند کارهای محاسباتی سنگین خود را بدون داشتن سخت افزار و نرم‌افزارهای گران، از طریق خدماتی انجام دهند. رایانش ابری آخرین پاسخ فناوری به این نیازها بوده است [۱].

تامین منابع برنامه‌های کاربردی چندلایه در محیط‌های ابری با یکسری چالش‌ها روبروست که شامل کسر تامین، اضافه تامین و نوسان است. در این مقاله برای مرتفع کردن چالش‌های مطرح‌شده همچنین بهینه‌سازی زمان‌بندی درخواست‌های کاربران و پاسخ‌دهی به آن‌ها، تقلیل مشکل تخطی از سرویس، به ارائه رویکردی بهبود یافته مبتنی بر یادگیری ماشین با بهره‌گیری از حلقه "MAPE"^۱ می‌پردازیم. در مرحله تحلیل این حلقه از مدل رگرسیون خطی (LRM) و در مرحله برنامه‌ریزی، از روش مبتنی بر نظریه بیز به منظور بهینه نمودن اقدامات استفاده شده است. سپس رویکرد پیشنهادی خود را تحت بار کاری واقعی FIFA با روش‌های Stat-RA و DPM-RA مقایسه نموده‌ایم که راهکار ارائه شده نسبت به راهکارهای پیشین، منجر به افزایش تعداد ماشین‌های مجازی به میزان ۱۰ درصد با بهبود نرخ مقیاس‌بندی، کاهش ۸ درصدی

2- Linear Regression
3- Bayes theorem

* نویسنده مسئول
1- Monitor, Analyze, Plan, Execute_Knowledge (MAPE-K)

رایانش ابری، باید به صورت خودکار منابع مورد نیاز به کاربران را تخصیص دهد و اگر کاربری مقداری بیشتر از منابع را خواستار باشد، باید آن منابع بدون آن که سرویس کاربر با مشکل مواجه شود به مشتری تخصیص یابد و در صورتی که منابع مورد استفاده بیش از حد مورد نیاز باشد باید منابع اضافی به صورت موقت خاموش شده تا در صورت نیاز مجدداً مورد استفاده قرار گیرند. این مفاهیم در دنیای رایانش ابری تحت عنوان «مقیاس پذیری»^۴ و یا گاهی اوقات «مقیاس بندی»^۵ شناخته می شود. تعاریف گوناگونی برای مقیاس پذیری وجود دارد که به چند مورد آن اشاره می کنیم [۲].

از آنجائی که برنامه های کاربردی، مخصوصاً برنامه های کاربردی تحت وب، دارای الگوهای حجم کاری منظمی نیستند، بنابراین عملیات مقیاس بندی (افزایش یا کاهش مقیاس) باید صورت بیدرنگ و با حداقل دخالت انسان انجام گیرد تا منابع در اسرع وقت برای برنامه های کاربردی تامین شود. به این نوع مقیاس بندی کردن منابع، که با کمترین دخالت انسان و به صورت خودکار انجام می شود، مقیاس بندی خودکار^۶ گفته می شود. مقیاس بندی خودکار، راهبرد اصلی در راستای تامین پویای منابع برای برنامه های کاربردی ابری است. فرایند مقیاس بندی، توسط عاملی به نام Auto-Scaler انجام می شود [۳]. Auto-Scaler، مسئول تصمیم گیری در مورد اعمال مقیاس بندی، بدون تعامل و مداخله یک مدیر انسانی است. هدف Auto-Scaler، تطبیق پویای منابع انتساب داده شده به برنامه های کاربردی بر اساس بارکاری ورودی است و باید توازنی بین برآورده کردن اهداف SLA و حداقل سازی هزینه، برقرار نماید. هر Auto-Scaler، با چالش های زیر روبرو خواهد شد:

کسر تامین^۷: در این مورد، برنامه کاربردی، منابع کافی برای پردازش همه درخواست های ورودی با توجه به

محدودیت های زمانی اعمال شده توسط SLA، را ندارد. در مورد ترافیک های ناگهانی، کسر تامین منابع برای برنامه کاربردی ممکن است منجر به تخطی از SLA شود.

اضافه تامین^۸: در این مورد، برنامه کاربردی منابع بیشتر از نیاز برای برآورده کردن توافقات SLA دارد. اگر چه این حالت، برای برآورده کردن اهداف SLA مناسب است اما به دلیل وجود منابع بیکار، مشتری باید هزینه های غیر ضروری پرداخت نماید.

نوسان^۹: این مورد، ترکیبی از دو پدیده نامطلوب کسر تامین و اضافه تامین است. در محیط ابری، وقتی با وجود اضافه تامین منابع برای یک برنامه کاربردی، تخطی از SLA رخ دهد، به این پدیده، نوسان گفته می شود (شکل ۱).

در این مقاله به ارایه رویکردی نوین برای تامین کارآمد منابع برنامه های کاربردی در محیط ابری پرداخته ایم. روش پیشنهادی برای تامین منابع بر روی ساختار MAPE که شامل چهار مرحله پایش، تحلیل، برنامه ریزی و اجرا است [۵] ارایه شده که در مرحله تحلیل از «مدل رگرسیون خطی (LRM)»^{۱۰} و در مرحله برنامه ریزی از «نظریه بیز»^{۱۱} استفاده کرده ایم.

سهام های عمده این تحقیق عبارتند از:

- ارائه چارچوبی مبتنی بر حلقه MAPE برای تامین خودکار منابع برای برنامه های کاربردی چند لایه
- ارائه فرمول بندی برای مسئله تامین خودکار منابع برای برنامه های کاربردی چند لایه
- شبیه سازی رویکرد پیشنهادی تحت بارکاری واقعی مسئله ای که در خور توجه است، تامین منابع ابری به صورت بهینه می باشد تا تمام کاربران بتوانند از هر منبعی بر اساس میزان مورد نیازشان استفاده کنند. در صورتی که این منابع به صورت بهینه برای کاربران تامین و مدیریت شوند شاهد تغییرات قابل توجهی در تعداد ماشین های مجازی، میانگین بهره وری، زمان پاسخ دهی،

8- Over-Provisioning
9- Oscillation
10- Linear regression model
11- Queueing theory

4- Scalability
5- Scaling
6- Auto-Scaling
7- Under-Provisioning

هزینه و سود حاصله خواهیم بود.

ادامه این مقاله بدین صورت سازماندهی شده است: بخش دوم به کارهای مربوطه اختصاص دارد، در بخش سوم رویکرد پیشنهادی را که شامل چارچوب پیشنهادی و الگوریتم پیشنهادی است را تشریح خواهیم نمود و در بخش چهارم، به شبیه‌سازی و ارزیابی کارایی روش پیشنهادی خواهیم پرداخت و نهایتاً بخش پنجم به نتیجه‌گیری و پیشنهادهای اختصاص دارد.

۲- کارهای مربوطه

تحقیقات متنوع و زیادی در سال‌های اخیر در رابطه با تامین منابع برای برنامه‌های کاربردی چندلایه صورت گرفته است، بیشتر کارهای موجود براساس تکنیک‌هایی است که در خصوص مدیریت منابع استفاده شده است که هر یک از روش‌های موجود مزایا و معایبی دارد که بر تصمیم‌گیری در مورد انتخاب تکنیک، تاثیر خواهد گذاشت. در این بخش مروری بر تحقیقات اخیر در زمینه تامین منابع برنامه‌های کاربردی چندلایه خواهیم داشت.

جینگ‌بی و همکارانش در سال ۲۰۱۳، به مسئله تخصیص منابع مجازی‌شده برای مراکز داده ابری در مقیاس بزرگ «CDC»^{۱۲}ها [۶] پرداختند. آن‌ها یک معماری مدیریت خودمختار را بر اساس سازوکار مجازی‌سازی برای برنامه‌های وب چندلایه در CDCها به منظور افزایش تخصیص منابع ارائه دادند که این معماری می‌تواند اجرای موثر محیط مجازی شده خدمات کاربردی «VASE»^{۱۳} و نیازمندی‌های کارخواه‌های SLA را تضمین کند سپس یک مدل صف ترکیبی انعطاف‌پذیر ارائه دادند و به‌طور کامل به زمان پاسخ، توان عملیاتی، هزینه و درآمد مشخص شده SLA پرداختند. پس از آن یک مسئله بهینه‌سازی غیرخطی محدود شده را بیان کردند.

از الگوریتم مکاشفه‌ای با هدف حداکثر کردن سود کل «CIP»^{۱۴}ها جهت حل مسئله استفاده کردند که این الگوریتم

بر اساس تغییرات پویای حجم کار منابع مجازی شده را تخصیص می‌دهد. آن‌ها تاثیر راهبرد تخصیص پویای منبع را بر اساس ماشین‌های مجازی «DVM-RA»^{۱۵} با دو راهبرد «Stat-RA»^{۱۶} و «DPM-RA»^{۱۷} مقایسه کردند.

تامین پویای منبع یک تکنیک چالشی برای نیازمندی‌های SLA از برنامه‌های کاربردی چندلایه در ابر مبتنی بر مجازی‌سازی است. در همین راستا Heng Wu و همکارانش در سال ۲۰۱۳ [۷] به ارائه یک روش آگاه از سود با نظریه کنترل بازخورد پرداختند که طبق روش ارائه شده، آن‌ها هزینه تامین منبع در مقایسه با روش هزینه آشکار تا ۳۰٪ کاهش می‌یابد و به‌طور موثر می‌تواند موارد نقض SLA را در مقایسه با روش آگاه از هزینه کاهش دهد.

تغییرات درخواست‌ها از سوی مستاجران و استفاده‌کنندگان متفاوت است بنابراین انعطاف‌پذیری برای آن‌ها افزایش می‌یابد. آن‌ها یک مدل تخریب کارایی را مبتنی بر یادگیری ماشین جهت محاسبه هزینه انتقال ایجاد کردند سپس یک تابع سود و هزینه را برای همه پیکربندی‌های مجدد ارائه کردند. در نتیجه حداقل هزینه با بیشترین سود را به‌دست آوردند.

Bhuvan Urgaonkar و همکاران [۸]، تکنیکی پویا برای برنامه‌های چندلایه اینترنت را ارائه نمودند که شامل یک مدل صف انعطاف‌پذیر برای تعیین چگونگی تخصیص منابع به هر لایه برنامه کاربردی و ترکیبی از روش‌های پیش‌بینی و واکنش‌پذیر جهت تعیین زمان به منظور تامین منابع در مقیاس‌های زمانی کوچک و بزرگ است. با توجه به روشی که آن‌ها ارائه دادند اهداف زمان پاسخ حفظ می‌شود و سربار کاهش می‌یابد همچنین سربار تعویض کارسازها از چندین دقیقه به کمتر از یک ثانیه می‌رسد.

باتوجه به این که تغییرات حجم کار پویاست و تغییرات بارکاری در زمان‌های مختلف متفاوت است، بنابراین تکنیک پیشنهادی آن‌ها مناسب است. همچنین سازوکارهای تامین

15- Dynamically Resource Allocation Strategy Based On Virtual Machine

16- Static Resource Allocation Strategy

17- Dynamically Resource Allocation Strategy Based On Physical Machine

12- Cloud Data Center

13- Virtualised Application Service Environment

14- Cloud Infrastructure Provider

پیش‌بینی و واکنشی به‌عنوان معماری مرکز داده بر اساس پایش ماشین مجازی ارائه شده است.

کیوان رحیم‌زاده و همکاران در سال ۲۰۱۵ [۹]، یک مدل تحلیلی بر اساس "QN"^{۱۸} به منظور تخمین معیارهای کیفیت سرویس برای برنامه‌های کاربردی چندلایه در یک مرکز داده مجازی ارائه نمودند. همچنین یک متدولوژی جهت اندازه‌گیری ویژگی‌های "VMTA"^{۱۹} از درخواست‌ها بین لایه‌ها و معیارهای عملکرد برنامه‌های کاربردی ارائه نمودند سپس آزمایش‌های مختلفی را به منظور ارزیابی عملکرد تلفیقی لایه‌های VMTAها، در خصوص حجم کار تدریجی و انفجاری انجام دادند.

نتایج کارهای آن‌ها نشان می‌دهد که قابلیت و ویژگی‌های لایه‌ها مستقیماً بر عملکرد کلی VMTAها تاثیر می‌گذارد و مغایرت درخواست‌ها برای داشتن منابع نرم‌افزار و سخت‌افزار موجب افزایش VM مسدود شده می‌گردد که سبب کاهش عملکرد VMTA خصوصاً، در زمان حجم کار بالا می‌گردد. استفاده از این تکنیک موجب کاهش هزینه، امنیت سرویس‌های مجازی و تضمین اهداف SLA گردیده است.

مارتا بلتران در سال ۲۰۱۵ [۱۰]، جهت تامین منابع برنامه‌های کاربردی چندلایه، تکنیکی را ارائه داده است که در آن تامین برنامه به‌صورت خودکار و مجازی صورت می‌پذیرد. تمام تصمیم‌گیری‌ها نیز به‌صورت کاملاً پویا و بر اساس شرایط و بدون دخالت انسان انجام می‌شود. تکنیکی که توسط وی ارائه گردیده است، AutoMap نام دارد.

این تکنیک توانایی تطابق با SLA در هزینه‌های اضافی را با توجه به سربارها دارد. از مقیاس‌گذاری افقی و عمودی نیز استفاده شده است. این تکنیک برای تعیین مقدار منابع مجازی مورد نیاز در هر لایه از برنامه کاربردی، مصرف منابع را به حداقل می‌رساند و توانایی تعامل با انواع برنامه‌های کاربردی مختلف را دارد، وی از

چارچوب CloudSim جهت پیاده‌سازی تکنیک پیشنهادی خود استفاده نمود.

Sukhpal Singh و همکاران در سال ۲۰۱۵ [۱۱]، به بررسی تامین منابع مناسب در حجم کارهای ابری وابسته به نیازمندی‌های QoS پرداختند. پارامترهای QoS بر اساس تکنیک تامین منبع به منظور تامین موثر منابع مورد نیاز است. آن‌ها به تحلیل و طبقه‌بندی بر اساس الگوهای رایج و تامین حجم کار ابر قبل از برنامه‌ریزی واقعی پرداختند. از آنجایی که پیچیدگی مدیریت منابع در ابرها روز به روز در حال افزایش است، یک تکنیک موثر به منظور مدیریت منابع نیاز است.

این تکنیک‌ها به شناسایی حجم کار ابر، تحلیل حجم کار ابر به منظور یافتن نیازمندی‌های QoS، طبقه‌بندی حجم کار و تامین منابع حجم کار بدون نقض SLA می‌پردازند. آن‌ها خوشه‌بندی حجم کار را از طریق k-means مبتنی بر الگوریتم خوشه‌بندی انجام دادند و از نرم‌افزار شبیه‌سازی CloudSim به منظور شبیه‌سازی استفاده کردند و در نهایت توانستند زمان صف‌بندی، زمان اجرا و هزینه را کاهش دهند و رضایت QoS و رضایت مشتری را تامین کنند.

رویکرد ارائه شده توسط سمانه کربلایی مهدی و مصطفی قبائی آرانسی در [۱۲] دارای سه مرحله پایش، تحلیل و اجراست که در مرحله تحلیل از سیستم استنتاج عصبی-مرحله‌ای تطبیق‌پذیر استفاده می‌شود. این سیستم از قدرت آموزش شبکه عصبی و مزیت زبانی سیستم‌های مرحله‌ای به منظور تحلیل فرآیندهای پیچیده به‌صورت بسیار قدرتمند عمل می‌کند. سپس رویکرد پیشنهادی را تحت بارهای کاری واقعی و مصنوعی با دو روش "ADRP_Fuzzy"^{۲۰} و "ADRP_AL"^{۲۱} مورد ارزیابی و مقایسه قرار داده‌اند. نتایج به‌دست آمده نشان می‌دهد که رویکرد پیشنهادی موجب کاهش زمان پاسخگویی و افزایش بهره‌وری می‌گردد.

20- Automatic Dynamic Resource Provisioning Fuzzy

21- Automatic Dynamic Resource Provisioning Automata Learning

18- Queuing Network

19-Virtualized Multi-Tier Application

ریحانه خرسند و همکاران در سال ۲۰۱۸ [۱۳] تکنیکی جهت تأمین منابع ترکیبی برای برنامه‌های چندتایی ارائه داده‌اند که مبتنی بر ترکیبی از مفهوم محاسبات خودمختار، مدل پیش‌بینی رگرسیون بردار پشتیبانی "SVR"^{۲۲} و رویکرد فرآیند سلسله‌مراتبی تحلیلی مرحله‌ی "FAHP"^{۲۳} است چارچوبی که ارائه شده بر اساس حلقه کنترل MAPE-k برای تأمین منابع خودمختار برنامه‌های کاربردی چندلایه در محیط‌های محاسبات ابری است. اثربخشی این رویکرد پیشنهادی تحت بار کار واقعی و مصنوعی مورد بررسی شده و نتایج تجربی نشان می‌دهد که راه حل پیشنهادی از لحاظ اصولی فراتر است از ماشین‌های مجازی اختصاص یافته، زمان پاسخ و هزینه در مقایسه با دیگر رویکردهاست. محمدصادق اصلانپور و همکاران در سال ۲۰۱۷ [۱۴] به ارائه مقیاس خودکار برنامه‌های وب در ابرها با توجه به هزینه رویکرد آگاهانه پرداختند که در این روش مجری حرفه‌ای اهمیت و اثربخشی این مرحله از چرخه کنترل را نشان می‌دهد. برخلاف معمول که جریان، راه حل پیشنهادی دستورات مقیاس پایین را از طریق انتخاب آگاهانه ماشین‌های مجازی مازاد اجرا می‌کند. علاوه بر این، با آن ویژگی‌های جدید، ماشین‌های مجازی مازاد در حداکثر مدت صورت حساب خود در قرنطینه نگهداشته‌اند تا هزینه به حداکثر برسد نتایج شبیه‌سازی نشان می‌دهد که مجری پیشنهادی ضمن افزایش کارایی، هزینه اجاره ماشین‌های مجازی را به میزان ۷ درصد کاهش داده است. جینگ‌جینگ گو و همکاران در سال ۲۰۱۹ [۱۵] روندی را به منظور تأمین منابع تقاضا بر اساس تخمین بار و هزینه خدمات در محیط ابری ارائه نموده‌اند در این روش تقاضا برای منابع باید از قبل تخمین زده شود. برای این منظور، یک تخمین بارمدل بر اساس مدل "ARIMA"^{۲۴} و شبکه عصبی "BP"^{۲۵} ارائه شده است. این مدل می‌تواند بار را تخمین بزند و با توجه به داده‌های پیشینه، خطای

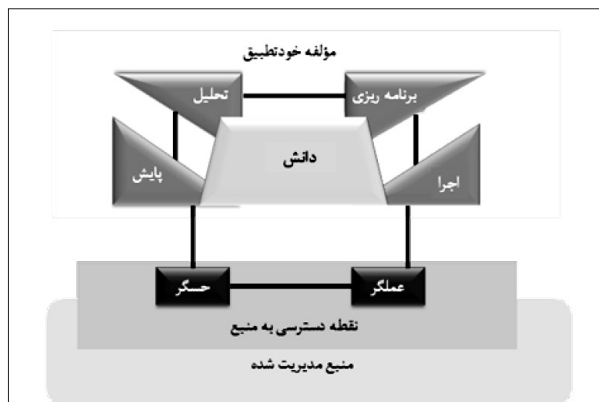
22- Support vector machines
 23- Fuzzy Analytic Hierarchy Process
 24- Auto Regressive Integrated Moving Average
 25- Back Propagation Network

تخمین را کاهش می‌دهد. قبل از انتشار منابع گره، کاربر برای اطمینان از عدم گم شدن اطلاعات کاربر، گره‌ها باید به گره‌های دیگر منتقل شوند. در این روش هنگام انتخاب هدف مهاجرت، سه معیار تعادل بار، زمان مهاجرت مصرف و هزینه مهاجرت خوشه در نظر گرفته شده است. با مقایسه نتایج تجربی، روش‌های پیشنهادی می‌توانند به‌طور موثر هزینه‌های خدمات را کاهش داده و خوشه را در حالت تعادل بار قرار دهد.

مطالعه تحقیقات فعلی نشان می‌دهد که کسر تامین، اضافه‌تأمین و نوسان در تأمین منابع برنامه‌های کاربردی چندلایه به‌عنوان مسئله‌ای چالش برانگیز مطرح بوده است زیرا منابع باید به‌صورت بهینه تأمین شوند تا تمام کاربران بتوانند از هر منبعی بر اساس میزان مورد نیازشان استفاده کنند. رویکردهای قبلی، هر کدام مزایا و معایبی دارند که آن‌ها را به‌طور مختصر بیان کرده‌ایم. بنابراین جهت بهینه‌سازی زمان‌بندی درخواست‌های کاربران و پاسخ‌دهی به آن‌ها، همچنین تقلیل مشکل تخطی از سرویس، در روش پیشنهادی این مقاله در مرحله تحلیل حلقه MAPE-k از مدل رگرسیون خطی (LRM) و در مرحله برنامه‌ریزی، از روش مبتنی بر نظریه بیز به منظور بهینه نمودن اقدامات استفاده شده است که بعد از ارزیابی و مقایسه با ارائه جداول و نمودارهای همه جانبه به تشریح قابلیت‌ها و بهینگی روش پیشنهادی می‌پردازیم و شاهد افزایش تعداد ماشین‌های مجازی با بهبود نرخ مقیاس‌بندی، کاهش میانگین بهره‌وری، کاهش زمان پاسخ‌دهی، در نتیجه کاهش هزینه تمام شده و افزایش سود حاصل خواهیم بود.

۳- دانش پیش زمینه

رایانش ابری، امکان دسترسی به محدوده وسیعی از منابع مجازی (مانند سخت‌افزار، نرم‌افزار، سرویس‌ها و...) را فراهم می‌آورد که این منابع می‌توانند به‌صورت پویا و با حجم متغییر به کاربران خود سرویس دهند، یعنی هر کاربر به اندازه‌ای که از آن منابع نیاز دارد، از آن‌ها استفاده



شکل ۱: چرخه MAPE-K [۱۶]

دخالت کاربر بتواند بر روی مدیریت معماری، عملکرد، کارایی و کنترل خرابی خود تصمیم گیری کند. در واقع در یک سیستم خودتطبیق معمولاً از سیستم‌های پشتیبانی تصمیم استفاده می‌شود. در یک سیستم خودتطبیق چرخه چهارگانه MAPE-K وجود دارد که عبارتند از پایش سیستم، تحلیل سیستم، برنامه‌ریزی و اجرا که همه این موارد در کنار یک پایگاه دانش انجام می‌شوند. این چرخه در شکل (۱). قابل مشاهده می‌باشد [۱۶]. این مدل بیشتر و بیشتر برای برقراری ارتباط بین جنبه‌های معماری سیستم‌های خودمختار استفاده می‌شود. به همین ترتیب این راهی روشن برای شناسایی و طبقه‌بندی بسیاری از کارهایی است که در این زمینه در حال انجام است.

بنابراین این بخش، چرخه MAPE-K را با جزئیات بیشتر معرفی می‌کند و پس از آن هر یک از اجزای سازنده آن را به نوبه خود در نظر می‌گیرد. ما بعد از آن همان کار را در نظر گرفته و از نقطه نظر قواعد خودمختاری به بررسی آن می‌پردازیم.

در چرخه MAPE-K یک عامل هوشمند، محیط خود از طریق حسگرها درک نموده و از این ادراک برای تعیین اقدامات به منظور اجرا بر روی محیط استفاده می‌کند.

در چرخه MAPE-K مؤلفه خودمختار که بلوک سازنده یک سیستم خودمختار است از مدیر خودمختار و مؤلفه مدیریت شده متشکل می‌باشد. عنصر مدیریت شده نشان دهنده هر منبع نرم‌افزاری یا سخت‌افزاری است

می‌کند که این کار به استفاده بهینه از منابع منجر می‌شود. رایانش ابری، باید به صورت خودکار منابع مورد نیاز به کاربران را تخصیص دهد و اگر کاربری مقداری بیشتر از منابع را خواستار باشد، باید آن منابع بدون آن که سرویس کاربر با مشکل مواجه شود به مشتری تخصیص یابد و در صورتی که منابع مورد استفاده بیش از حد مورد نیاز باشد باید منابع اضافی به صورت موقت خاموش شده تا در صورت نیاز مجدداً مورد استفاده قرار گیرند. این مفاهیم در رایانش ابری تحت عنوان مقیاس‌پذیری و یا گاهی اوقات مقیاس‌بندی شناخته می‌شود.

مقیاس‌پذیری یکی از بهترین راه‌های افزایش بهره‌وری و بهبود عملکرد سیستم‌های ابری می‌باشد و توسعه‌دهندگان ابر معتقدند که باید برای کاربران به مقدار نیاز منابع فراهم شود و از طرفی، کاربران نیز علاقه‌مند به پرداخت هزینه به مقدار استفاده شده (پرداخت-به‌ازای-استفاده) هستند. باید به این نکته توجه داشت که اکثر عناصری که در حوزه فناوری اطلاعات مطرح هستند، امکان تغییر مقیاس را دارا می‌باشند اما برخی از مهم‌ترین این عناصر عبارتند از: مقدار پردازنده اختصاص یافته، مقدار حافظه RAM در دسترس، مقدار فضای دیسک سخت، تعداد برنامه‌های اجرا شده، مقدار پهنای باند ارتباطی، تعداد کاربران نهایی، تعداد درگاه‌های ارتباطی، تعداد مجوزهای نرم‌افزاری و غیره. سازوکار تامین پویای منابع در رایانش ابری، به برنامه‌ها این امکان را می‌دهد تا در مقیاس‌های بالا و پایین اجرا شوند و بدین ترتیب عملکرد و پایداری اقتصادی آن‌ها متعادل می‌گردد. یکی از چالش‌های مهم در رایانش ابری، دسترس‌پذیری و مدیریت منابع می‌باشد. نیاز هست که سیستم‌های ابری، قوی و قابل دسترس باشند تا مدیریت به راحتی صورت گیرد. چرا که بیشتر تجارت‌ها، توسط شرکت‌های مختلف وابسته به خدمات ارائه شده است [۲].

۳-۱- چرخه خودمختار MAPE-K

یک سیستم خودتطبیق سیستمی است که بدون

که رفتار خودمختاری را از طریق قرار گرفتن در کنار مدیر خودمختار ارائه می‌دهد. بنابراین، عنصر مدیریت شده می‌تواند برای مثال یک کارساز وب و یا پایگاه داده، یکی از اجزای نرم‌افزاری خاص در یک برنامه (به‌عنوان مثال، بهینه‌ساز پرس‌وجو در پایگاه داده)، سیستم عامل، یک خوشه از ماشین‌ها در یک محیط شبکه، یک پشته از دیسک‌های سخت، یک شبکه سیمی یا بی‌سیم، CPU، چاپگر، و غیره باشد. مدیر خودمختار نیز مؤلفه نرم‌افزاری است که توسط مدیران و با استفاده از اهداف سطح بالا پیکربندی شده است، در هر دو مدل، کارساز، ابزاری برای جمع‌آوری اطلاعات مرتبط و تأثیرگذار بوده که برای یک سرویس‌دهنده وب، این می‌تواند شامل زمان پاسخ به درخواست مشتری، استفاده از شبکه و مصرف دیسک، CPU و مصرف حافظه باشد همچنین مسئول اعمال تغییراتی از قبیل اضافه کردن یا حذف کارساز به/ از یک خوشه کارساز وب بزرگ بوده یا شامل پارامترهای پیکربندی در یک کارساز وب کوچک در مؤلفه مورد نظر باشد.

چهار مرحله اصلی رایانش خودمختار پایش، تحلیل، برنامه‌ریزی و اجراست. که در ادامه مقاله آن‌ها را تشریح می‌کنیم.

۳-۱-۱ پایش

پایش شامل اتخاذ خواص محیط (فیزیکی یا مجازی، به‌عنوان مثال، یک شبکه)، که برای چهار خواص سیستم ذکر شده (خود-X) دارای اهمیت هستند می‌باشد. مؤلفه نرم‌افزاری و یا سخت‌افزاری مورد استفاده برای انجام پایش، حسگر نام دارد. به‌عنوان مثال، زمان تأخیر شبکه و پهنای باند عملکرد کارساز وب را اندازه‌گیری می‌کند، در حالی که نمایه‌سازی پایگاه داده و بهینه‌سازی پرس‌وجو زمان پاسخ "DBMS"^{۲۷} را تحت تأثیر قرار می‌دهند، که می‌توان آن‌ها را تحت پایش قرار داد. مدیر خودمختار برای تشخیص خرابی یا عملکرد بهینه عنصر خودمختار به اطلاعات پایش شده مناسب و تغییرات مناسب تأثیرگذار

27- Data-base Management System

نیاز دارد. انواع خواص پایش شده، و حسگرهای مورد استفاده، اغلب برنامه کاربردی خاصی هستند، و عوامل مؤثر مورد استفاده برای اجرای تغییرات در عنصر مدیریت شده نیز برنامه کاربردی خاص هستند.

۳-۱-۲ تحلیل

مرحله تحلیل، با پردازش معیارهای جمع‌آوری شده از سیستم پایش، سر و کار دارد و با پردازش این معیارها، داده‌هایی در مورد وضعیت بهره‌وری فعلی سیستم و پیش‌بینی‌هایی از نیازهای آتی، به‌دست می‌آورد تا در صورت نیاز به مختار، پاسخ مناسبی به تحریکات داده شود. بعضی از رایانش‌های خودمختار هیچ نوع پیش‌بینی انجام نمی‌دهند و فقط به وضعیت جاری سیستم پاسخ می‌دهند که به آن‌ها واکنشی می‌گوییم، اما برخی دیگر از این محاسبات از تکنیک‌های پیچیده‌ای برای پیش‌بینی تقاضاهای آتی به منظور تنظیم منابع به اندازه پیش‌بینی شده، استفاده می‌کنند که به آن‌ها پیش‌کنشی می‌گوییم. پیش‌بینی مهم است زیرا همیشه یک تاخیر زمانی از لحظه اجرا شدن عمل رایانش خودمختار (مثلاً اضافه کردن ماشین مجازی) تا مؤثر واقع شدن عمل (به دلیل راه‌اندازی مجدد سیستم، برای عملیاتی شدن یک ماشین مجازی چند دقیقه زمان لازم است) وجود دارد. سیستم‌های واکنشی، قادر به مقیاس‌بندی ترافیک‌های ناگهانی بار کاری ورودی نیستند و به همین دلیل از رویکردهای پیش‌کنشی، به منظور برخورد با نوسانات تقاضا و تامین پیشاپیش منابع استفاده می‌شود.

۳-۱-۳ برنامه‌ریزی

این مرحله، هسته اصلی رایانش خودمختار است که در وسیع‌ترین مفهوم، شامل اتخاذ داده‌های پایش از حسگرها برای تولید یک سری از تغییرات برای تأثیر بر عنصر مدیریت شده می‌باشد تا در جهت بهبود عملیات ارائه شده، تغییراتی اعمال شود.

اگر اعمال این روش به شیوه‌ای بدون حالت (که در آن مدیر خودمختار هیچ اطلاعاتی در مورد حالت عنصر

مدیریت شده را نگه نمی‌دارد و برای تصمیم‌گیری مبنی بر این که آیا بر یک طرح خودمختاری اثر گذارند یا خیر، صرفاً متکی بر اطلاعات حسگر جریان می‌باشد) بسیار محدود است. در واقع، برای مدیر خودمختار به مراتب بهتر است که اطلاعات را بر روی حالت عنصر مدیریت شده که می‌تواند به تدریج از طریق داده‌های حسگر به روز شود، نگه دارد تا با توجه به وضعیت سیستم و سایر اطلاعات موجود، تصمیم بگیرد چه راهبرد بایستی اتخاذ شود.

با توجه بیشتر به موضوع اطلاعات حالتی که مدیر خودمختار باید در مورد عنصر مدیریت شده حفظ کند، تحقیقات زیادی یک رویکرد کامل‌تر برای مدیریت یک سیستم (که رویکرد مدل‌محور معماری نام دارد) که در آن نوعی از مدل کل سیستم مدیریت شده توسط مدیر خودمختار ایجاد می‌شود، را مورد بررسی قرار دادند. این معمولاً یک مدل معماری نامیده می‌شود که نشان دهنده رفتار سیستم مدیریت شده، الزامات آن، و احتمالاً هدف آن می‌باشد. مدل همچنین ممکن است نشان دهنده برخی از جنبه‌های محیط عملیاتی باشد که در آن عناصر مدیریت شده قرار دارند که در آن محیط عملیاتی می‌تواند به‌عنوان هر خاصیت قابل مشاهده (توسط حسگرها) درک شود که به‌عنوان مثال می‌تواند اجرا، ورودی کاربر نهایی، دستگاه‌های سخت افزاری و خصوصیات اتصال به شبکه را تحت تاثیر قرار دهد.

مدل از طریق داده‌های حسگر به‌روز شده و به همین دلیل در سیستم مدیریت شده برای برنامه‌ریزی سازگاری مورد استفاده قرار می‌گیرد. یک مزیت بزرگ رویکرد مبتنی بر مدل معماری برای برنامه‌ریزی این است که با این فرض که مدل به درستی منعکس کننده سیستم مدیریت شده باشد، مدل معماری را می‌توان برای تأیید این که یکپارچگی سیستم هنگام اعمال یک سازگاری حفظ می‌شود، استفاده نمود. به این معنا که می‌توانیم تضمین کنیم که پس از آن که سازگاری برنامه‌ریزی شده اجرا شود، سیستم به درستی به کار خود ادامه خواهد داد. این به‌خاطر این است

که تغییرات برنامه‌ریزی شده اول به مدل اعمال می‌شود، که حالت سیستم ناشی از سازگاری (از جمله هر گونه نقض محدودیت و یا الزامات سیستم ارائه شده در مدل) را نشان خواهد داد. اگر حالت جدید سیستم قابل قبول باشد، این طرح پس از آن می‌تواند بر روی سیستم مدیریت شده واقعی تاثیر گذارد، بنابراین اطمینان می‌دهد که مدل و پیاده سازی نسبت به یکدیگر سازگار هستند.

هنگامی که یک تغییر در سیستم مدیریت شده رخ می‌دهد و این تغییر به مدل اعمال می‌شود، همیشه یک تاخیر بین زمان این دو کار وجود دارد. در واقع، در صورتی که تاخیر به اندازه کافی بالا بوده و سیستم نیز مکرر تغییر نماید، ممکن است یک طرح سازگاری ایجاد شده و برای اجرا ارسال گردد.

۳-۱-۴ اجرا

این فرایند مسئول اجرای تصمیمات اتخاذ شده توسط مرحله برنامه‌ریزی است. به عبارتی دیگر، این مرحله، مرحله‌ای سر راست و مستقیم است که از طریق API ارائه‌دهنده ابری، پیاده‌سازی می‌شود و پیچیدگی‌های ذاتی و واقعی را از دید مشتری مخفی می‌کند [۵].

۳-۲ یادگیری بیزی

استدلال بیز روشی احتمالی برای استنتاج ارائه می‌کند. این روش بر اساس این فرض است که کمیت‌های مورد نظر از توزیع‌های احتمال پیروی می‌کنند و تصمیم‌گیری بهینه را می‌توان با استدلال بر این توزیع‌های احتمال و داده‌های مشاهده شده انجام داد. اهمیت این روش در یادگیری ماشین این است که روشی کمی برای ارزیابی مدارک فرضیه‌ها ارائه می‌کند. استدلال بیزی اساس الگوریتم‌های یادگیری که با استفاده از احتمالات کار می‌کنند است. همچنین استدلال بیز چارچوبی برای بررسی عملیات دیگر الگوریتم‌هایی که از احتمالات استفاده نمی‌کنند ایجاد می‌کند.

روش‌های یادگیری بیزی به دو دلیل به مطالعه ما در مورد یادگیری ماشین مربوط می‌شود. اول این که

الگوریتم های یادگیری بیزی احتمال صریح هر فرضیه، مثل دسته‌بندی‌کننده ساده بیز را محاسبه می‌کنند. این نوع روش‌ها از پرکاربردترین روش‌ها در بعضی مسائل یادگیری هستند [۱۷].

۳-۲-۱ ویژگی‌های یادگیری بیزی

آشنایی پایه‌ای با روش‌های بیزی برای درک بسیاری از الگوریتم‌های یادگیری ماشین اهمیت بسزایی دارد. ویژگی‌های روش‌های یادگیری بیزی شامل این موارد می‌باشد:

- هر نمونه آموزشی می‌تواند احتمال تخمینی این که فرضیه درست است را کم یا زیاد کند. این حقیقت باعث می‌شود که روش‌های بیزی نسبت به الگوریتم‌هایی که کاملاً فرضیه را با نمونه‌های غیرسازگار رد می‌کنند انعطاف‌پذیرتر باشند.

- می‌توان از دانش قبلی به همراه داده‌ای مشاهده شده برای تعیین احتمال نهایی درستی فرضیه‌ها استفاده کرد. در یادگیری بیزی، دانش قبلی با (۱) احتمال اولیه هر فرضیه و (۲) احتمال روی داده‌های تعیین شده برای هر فرضیه ممکن تعیین می‌شود.

- روش‌های بیزی می‌توانند برای فرضیه‌ها احتمالاتی را پیش‌بینی کنند (برای مثال، فرضیه این بیمار ذات‌الریه با احتمال ۹۳٪ کاملاً بهبود خواهد یافت).

- نمونه‌های جدید را می‌توان با ترکیب پیش‌بینی‌های چندین فرضیه، (هر کدام با وزن احتمال‌شان) دسته‌بندی کرد.

حتی هنگامی که اثبات می‌شود که روش‌های بیزی محاسباتی غیرقابل پیش‌بینی انجام می‌دهند، با این حال معیار استاندارد برای دیگر روش‌های عملی یادگیری مطرح می‌کنند [۱۸].

۳-۳ برنامه‌های کاربردی چند لایه

اخیراً پیشرفت‌های بزرگی در معماری سازمانی تحت وب صورت گرفته است. بخصوص در رایانش ابری که ارائه‌دهندگان خدمات را قادر می‌کند خدمات مبتنی بر

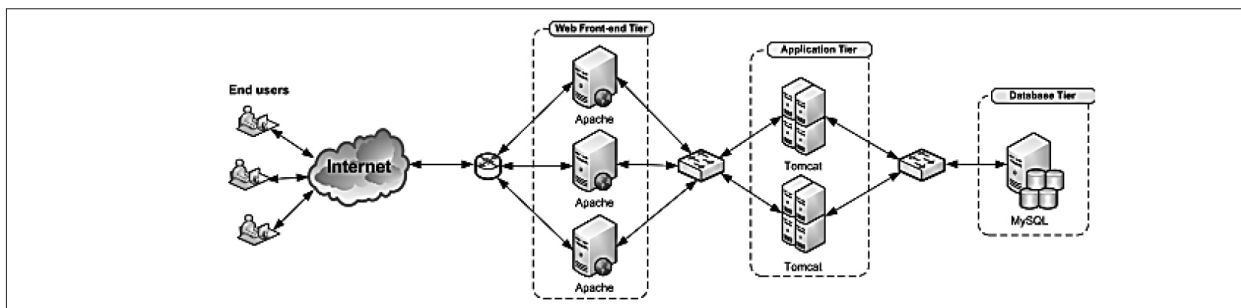
معماری وب سنتی را متوقف کنند. ترتیب تکامل چنین برنامه‌های کاربردی (از معماری رایانه بزرگ تا کارخواه/ کارساز سنتی و سپس پیاده‌سازی مجازی فعلی) گویای این است که راه‌اندازی سیستم مستقل نه تنها هزینه‌های کلی را افزایش می‌دهد بلکه کاهش مقیاس‌پذیری، دسترس‌پذیری و کاهش تحمل خطای سرویس را به همراه خواهد داشت. امروزه ارائه‌دهندگان خدمات عمدتاً سرویس‌های مبتنی بر وب، برنامه‌های کاربردی قدیمی و سیستم‌های پایگاه داده را روی مجموعه‌ای از کارسازها مستقر می‌کنند.

معماری تک‌لایه ساختار نسبتاً ساده‌ای دارد و راه‌اندازی آن آسان است، اما اغلب وبگاه‌های امروزی از معماری چندلایه استفاده می‌کنند. این معماری، پردازش برنامه کاربردی را به چندین لایه تقسیم می‌کند. هر لایه، قابلیت‌های مشخصی را ارائه می‌کند. مزیت چنین معماری این است که می‌تواند سطح بالایی از مقیاس‌پذیری و قابلیت اطمینان را فراهم کند [۱۹]. علاوه بر این، معماری چندلایه برای بهینه‌سازی هزینه‌های مالی در فرایند گردش کار نیز مفید است. با این وجود، به دلیل وابستگی لایه‌ها، تخصیص منابع میان آن‌ها مشکل است.

بنابراین ارائه‌دهندگان خدمات برای وبگاه‌های تجارت الکترونیکی مثل eBay و سیستم‌های اطلاعاتی مبتنی بر وب مانند ویکی‌پدیا، از معماری چند لایه استفاده می‌کنند.

معماری چند لایه شامل لایه‌های مختلف است که هر لایه میزبان کارسازهایی با عملکرد مشابه است. هر لایه یک سرویس از پیش تعریف شده برای جایگزینی و دریافت خدمات از لایه قبلی ارائه می‌کند. واضح است که اولین و آخرین لایه، قبل و جایگزین ندارند.

اغلب برنامه‌های کاربردی چند لایه تحت وب، از یک معماری سه‌لایه استفاده می‌کنند. سه لایه شامل لایه نمایش، لایه نرم‌افزار (منطق تجاری) و لایه دسترسی به داده به ترتیب به صورت کارساز وب، کارساز نرم‌افزار (برنامه کاربردی) و کارساز پایگاه داده، پیاده‌سازی شده‌اند.



شکل ۲: معماری برنامه کاربردی ۳- لایه [۹]

کارسازهای پایگاه داده است. این لایه، پردازش پایگاه داده و دسترسی به داده‌ها را مدیریت می‌کند. لایه کارساز پایگاه داده، برای ذخیره و بازیابی اطلاعات یک وبگاه (مانند حساب کاربری، فهرست گزارش‌ها و سفارش‌های مشتری) استفاده می‌شود.

بر اساس قابلیت‌های لایه‌ها، تقاضای منابع از هر لایه به لایه دیگر متفاوت است. به‌عنوان مثال لایه وب روی CPU و لایه پایگاه داده روی I/O متمرکز است. هر لایه نرم‌افزارها و سخت‌افزارهای مختلفی دارد و قادر است تعداد مشخصی از درخواست‌ها را سرویس‌دهی کند. با توجه به قابلیت‌های مختلف و تقاضا برای منابع، برنامه‌های کاربردی چند لایه با چالش‌هایی مانند وابستگی عملکرد و جابه‌جایی تنگنا بین لایه‌ها مواجه است. ارسال درخواست بیشتر در یک لایه ممکن است آن را مستعد تنگنا و مشکلات عملکرد کند. تخریب عملکرد یک لایه ممکن است باعث ایجاد تنگنا شود که عملکرد برنامه را به‌طور کلی مورد تهدید قرار می‌دهد. (مثل اثر دومینو). در چنین شرایطی برخی لایه‌ها ممکن است استفاده از منابع لایه‌های بالا را تجربه کنند؛ در حالی که لایه‌های دیگر مورد استفاده قرار گرفته‌اند و در انتظار دریافت درخواست لایه‌های دیگرند. در نهایت عملکرد نرم‌افزار به‌طور کلی تخریب خواهد شد. بنابراین طراحی یک برنامه کاربردی متعادل‌کننده بار مناسب می‌تواند از تخریب بیشتر جلوگیری کند [۹]. در این بخش قصد داریم رویکرد پیشنهادی خود را با جزئیات بیشتر توضیح دهیم، روشی برای تامین منابع بر روی ساختار MAPE-K با استفاده از ترکیب مدل رگرسیون خطی (LRM) و نظریه بیز

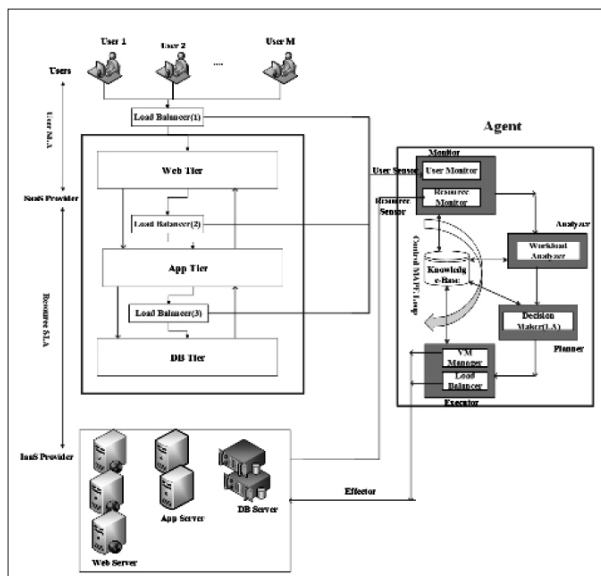
شکل (۲) یک برنامه کاربردی با یک پیکربندی دلخواه از کارسازها را در یک مرکز داده نشان می‌دهد. زنجیره‌ای از لایه‌ها ممکن است شامل متعادل‌کننده بار، پروکسی کارسازها، دروازه‌ها یا کارسازهای حافظه نهمان برای بهبود عملکرد نرم‌افزار و امنیت باشند. این مؤلفه‌ها در شکل به‌صورت اختصاری نشان داده شده‌اند.

اولین لایه که لایه نمایش نامیده می‌شود، شامل کارسازهای وب است. این لایه هر چیزی را که در سمت مشتری، داخل مرورگرهای وب ارائه می‌شود، نمایش می‌دهد. برای کارساز وب، سه وظیفه اصلی وجود دارد [۲۰]: (۱) پذیرفتن/ رد کردن درخواست‌های مشتریان و درخواست‌های سرویس‌های ایستای تحت وب؛ (۲) عبور دادن درخواست‌ها به سمت کارساز نرم‌افزار؛ (۳) دریافت پاسخ از کارساز نرم‌افزار و فرستادن آن‌ها برای مشتریان. کارساز آپاچی^{۲۸} و کارساز اطلاعات اینترنت مایکروسافت (IIS) نمونه‌هایی از این نوع کارسازها هستند.

دومین لایه که لایه نرم‌افزار نامیده می‌شود، شامل کارسازهای نرم‌افزار است. پردازش منطق تجاری، در این لایه اجرا می‌شود. همچنین سه وظیفه این لایه عبارت است از: (۱) دریافت درخواست‌ها از کارساز وب؛ (۲) جستجوی اطلاعات در پایگاه داده و پردازش اطلاعات؛ (۳) فرستادن اطلاعات پردازش شده به کارساز وب. کارسازهای نرم‌افزار اغلب از آپاچی تامکت^{۲۹} و کارساز نرم‌افزار سیستم سان‌جاوا^{۳۰} استفاده می‌کنند.

آخرین لایه که لایه داده نامیده می‌شود، شامل

28-Apache Server
29-Apache Tomcat
30-Sun Java



شکل ۳: چارچوب تأمین منابع برای برنامه‌های کاربردی چندلایه بر اساس حلقه کنترلی MAPE

تأمین منابع برای برنامه کاربردی چند لایه تأمین کننده SaaS بر اساس حلقه کنترلی MAPE-K می‌باشد که از طریق دو زیرمولفه پایش کاربر و پایش منبع درخواست‌های کاربران و وضعیت منابع را به ترتیب پایش کرده و از طریق دو زیرمولفه متوازن کننده بار و مدیر ماشین مجازی میزان منابع مناسب را به درخواست‌ها تخصیص می‌دهد. علاوه بر این، دو نوع SLA به نام‌های SLA کاربر و SLA منبع وجود دارد. SLA کاربر، قراردادی بین تأمین کننده SaaS و کاربر نهایی است و شامل فقره‌هایی مانند: مهلت اتمام، بودجه، نرخ جریمه و طول درخواست می‌باشد. SLA منبع، قراردادی بین تأمین کننده SaaS و تأمین کننده IaaS می‌باشد و شامل فقره‌هایی مانند: قیمت ماشین مجازی، نوع ماشین مجازی، سرعت پردازش، زمان سرویس و سرعت انتقال داده می‌باشد. در این پژوهش، برای محاسبه تخطی از SLA بر روی تعامل بین کاربر و تأمین کننده (SaaS (SLA) کاربر متمرکز می‌شویم.

همان‌طور که در شکل (۳) نشان داده شده است، مولفه‌های اصلی چارچوب تأمین منابع مبتنی بر حلقه کنترل MAPE-K عبارتند از:

● پایش^{۳۴}: مولفه پایش، مسئول جمع‌آوری معیارهای

با توجه به این که این ساختار شامل چهار مرحله پایش، تحلیل، برنامه‌ریزی و اجرا است، در روش پیشنهادی تمامی مرحله‌ها عملکرد خاصی دارند. در ادامه این بخش ابتدا چارچوب مبتنی بر این حلقه را شرح می‌دهیم و در نهایت الگوریتم پیشنهادی در این خصوص را ارائه می‌دهیم.

۴-۱ چهارچوب پیشنهادی

در این بخش، یک چهارچوب تأمین منبع برای برنامه‌های کاربردی چند لایه بر اساس حلقه کنترلی MAPE-K، برای پیاده‌سازی رویکرد پیشنهادی ارائه خواهیم داد. همان‌طور که در شکل (۲) نشان داده شده است، چهارچوب کلی شامل سه مولفه اصلی به نام تأمین کننده "SaaS"^{۳۱}، تأمین کننده "IaaS"^{۳۲} و مولفه عامل خودمختار مبتنی بر حلقه کنترلی MAPE-K می‌باشد. تأمین کننده SaaS به‌عنوان مالک برنامه کاربردی سه لایه برای ارائه سرویس به کاربران نهایی می‌باشد از آنجائی که ممکن است تعداد درخواست‌های ورودی به یک لایه از ظرفیت میزان ماشین‌های مجازی تأمین شده برای آن لایه، بیشتر باشد، صفی از درخواست‌ها قبل از ورود به هر لایه در قالب متوازن کننده بار^{۳۳} تشکیل می‌شود. کاربران نهایی، درخواست‌هایشان را برای استفاده از برنامه کاربردی ابری، به تأمین کننده SaaS ارائه می‌کنند. تأمین کننده SaaS، برای میزبانی برنامه کاربردی‌اش، یا از منابع مراکز داده داخلی خودش استفاده می‌کند یا این که منابع مورد نیاز برای میزبانی برنامه کاربردی را از یک تأمین کننده IaaS مانند Amazon EC2 اجاره می‌کند (که فرض ما در این تحقیق بر همین اساس است). تأمین کنندگان IaaS، منابع نامحدودی در قالب ماشین‌های مجازی (کارساز) در رده‌های مختلف (مثلاً سه نوع ماشین مجازی (Small, Medium, Large) برای اجرای درخواست‌های لایه‌های مختلف برنامه کاربردی چند لایه به تأمین کنندگان SaaS ارائه می‌دهد. مولفه عامل خودمختار مسئول مدیریت

31- Software as a Service
32- Infrastructure as a Service
33- Load Balancer

34- Monitor

مرتبط با منابع و کاربران است. مولفه پایش شامل دو زیرمولفه پایش کاربر و پایش منبع می‌باشد. مولفه پایش منبع، مسئول جمع‌آوری اطلاعات در مورد منابع محاسباتی، ذخیره‌سازی و شبکه (مانند بهره‌وری پردازنده، میزان حافظه مصرف‌شده، متوسط زمان سرویس و ترافیک شبکه) از تامین کننده IaaS می‌باشد. مولفه پایش کاربر، مسئول جمع‌آوری اطلاعات در مورد بارکاری و اگذار شده توسط کاربران (مانند نرخ درخواست‌های ورودی، نوع درخواست‌ها، اندازه درخواست‌ها، درخواست‌های پردازش شده، درخواست‌های متوقف شده و...) از سه صف درخواست (Load Balancer1, Load Balancer3, Load Balancer2) می‌باشد. این اطلاعات پایش شده توسط دو زیرمولفه پایش کاربر و پایش منبع، جمع‌آوری و یکپارچه شده و برای استفاده توسط سایر مراحل در پایگاه دانش ذخیره می‌شود.

● **تحلیل‌گر**^{۳۵}: مولفه تحلیل‌گر، مسئول پردازش اطلاعات جمع‌آوری شده از مولفه پایش می‌باشد. داده‌های به دست آمده توسط زیرمولفه‌های پایش، توسط مولفه تحلیل‌گر برای بررسی اعمال خودمختاری مورد نیاز برای تضمین سطح کیفیت سرویس، بررسی می‌شوند. مولفه تحلیل‌گر بارکاری، از تکنیک‌های پیش‌بینی برای پیش‌بینی تعداد درخواست‌ها و یا بهره‌وری پردازنده برای مقابله با نوسانات بارکاری استفاده می‌کند. علاوه بر این، تحلیل‌گر زمان پاسخ مشتری‌های مختلف را نیز دریافت می‌کند.

● **برنامه‌ریز**^{۳۶}: مولفه برنامه‌ریز، هسته اصلی چارچوب تامین منابع است. این مولفه تعیین می‌کند که چه موقع و چه تعداد ماشین مجازی باید به سرویس‌های ابری مختلف برای پیدا کردن یک مصالحه بین میزان هزینه و میزان تخطی از SLA تخصیص یابد. در این مقاله از تکنیک یادگیری بی‌یاز به‌عنوان تصمیم‌گیرنده در این مولفه استفاده می‌کنیم که از نتایج پیش‌بینی مولفه تحلیل‌گر برای حذف یا اضافه کردن ماشین مجازی استفاده می‌کنند.

● **اجراکننده**^{۳۷}: مولفه اجراکننده شامل دو زیرمولفه متوازن کننده بار و مدیر ماشین مجازی است. زیرمولفه متوازن کننده بار، تمام درخواست‌های ورودی از کاربران را دریافت کرده و آن‌ها را بین ماشین‌های مجازی مناسب بر طبق سیاست‌های توازن بار (مانند Round Robin, Random, ...) توزیع می‌کند. به دلیل این که زیرمولفه مدیر ماشین مجازی، به‌طور مستقیم با ماشین‌های مجازی تامین کنندگان IaaS در ارتباط است، این مولفه، مسئول اجرای واقعی اعمال (Scale up/Scale down) تصمیم گرفته شده توسط مولفه برنامه‌ریز می‌باشد.

۴-۲- الگوریتم پیشنهادی

در این بخش، الگوریتم تامین خودکار منبع برای لایه‌های مختلف برنامه کاربردی ابری را با جزئیات بیشتری تشریح می‌کنیم و به بیان الگوریتم اصلی می‌پردازیم که چهار تابع MAPE-K را فراخوانی می‌کند. پس از استقرار ارائه دهنده SaaS با لایه‌های مختلف ابر در چهارچوب ابر، الگوریتم ارائه شده تا زمانی که هیچ یک از لایه‌های ابر در ارائه دهنده SaaS در حال اجرا نباشد، اجرا می‌شود. این الگوریتم حلقه کنترل MAPE-K را دنبال می‌کند و ماشین‌های مجازی را مدیریت می‌نماید. این الگوریتم، ابتدا تعدادی از ماشین‌های مجازی مناسب را بر اساس تاریخچه تجربی موجود با استفاده از پیش‌بینی تکنیک‌ها برای هر لایه از مخزن VMها توسط ارائه دهنده IaaS راه‌اندازی می‌کند. برای هر لایه ابر Li توسط ارائه دهنده SaaS در Δt پیشنهاد شده است. حلقه کنترل MAPE-K همانند آنچه در الگوریتم (۱) نشان داده شده است به‌طور منظم در بازه‌های زمانی خاص اجرا می‌شود و با مقدار بهینه منابع بر اساس حجم کار لایه‌های ابر منطبق می‌گردد. در این الگوریتم حلقه کنترل MAPE-K شامل مراحل پایش، تحلیل، برنامه‌ریزی و اجرا می‌باشد. ابتدا پایش سیستم شروع به جمع‌آوری اطلاعات درباره کاربر و وضعیت منابع می‌کند سپس تحلیل‌گر از اطلاعات به‌دست آمده به منظور تخمین منبع

35- Analyzer
36- Planner

اضافی و نیازمندی‌ها در این خصوص استفاده می‌نماید. سپس در مرحله برنامه‌ریزی، اقدامات لازم در خصوص منابع صورت می‌گیرد یعنی در صورت نیاز ماشین مجازی اضافه یا حذف می‌شود. در نهایت، ارائه دهنده اقدامات درخواست شده توسط برنامه‌ریز را برای مرحله بعد اجرا می‌کند.

Algorithm1: Pseudocode for Automatic resource provisioning	
1:	Initialization: boots an appropriate number of VMs for cloud layers
2:	while (the system is running and in the beginning of interval Δt) do
3:	Begin
4:	for (every cloud layer L_i in multi layer cloud application at the interval Δt) do
5:	Begin
6:	Monitoring (L_i)
7:	Analysis (L_i)
8:	Planning (L_i)
9:	Execution ()
10:	end for
11:	end while

۴-۲-۱ مرحله پایش

مرحله پایش، مرحله‌ای است که معیارهای مرتبط با کاربران و منابع، پایش شده‌اند. الگوریتم (۲) به منظور عملیات پایش در پژوهش مورد نظر مورد استفاده قرار گرفته است. در بخش کاربر، برای هر لایه ابر L_i که توسط یک ارائه‌دهنده SaaS پیشنهاد شده است، پایش کاربر به معیارهایی از تعداد VMها که باید از ارائه‌دهنده IaaS اجاره شوند و توسط ارائه دهنده SaaS برای لایه ابر L_i در بازه Δt ام و تعداد درخواست‌ها برای لایه L_i در بازه Δt ام اجرا شود را پایش می‌کند. در بخش منابع، پایش منبع، مصرف پردازنده VMهایی که به لایه L_i در Δt تخصیص یافته‌اند، تعداد VMها و بارکاری را در Δt پایش می‌کند. این اطلاعات برای استفاده در سایر مراحل در پایگاه دانش ذخیره می‌شود.

Algorithm2: Pseudocode for Monitoring phase	
1:	Begin
2:	for (every cloud layer L_i in multi layer cloud application at the interval Δt) do
3:	Begin
4:	Monitor($W_i(\Delta t)$); /* workload monitoring*/
5:	Monitor($Num_i(\Delta t)$); /*Resource Monitoring*/
6:	Monitor($U_i(\Delta t)$); /*Monitoring CPU utilization of Layer L_i */
7:	end for
8:	end

در این مرحله، معیارهای به‌دست آمده از مرحله پایش مستقیماً به منظور پیش‌بینی میزان بهره‌وری پردازنده برای لایه ابر L_i در بازه $\Delta(t+1)$ با استفاده از تکنیک‌های پیش‌بینی پردازش می‌شوند. بنابراین تحلیلی را انجام می‌دهیم که بر اساس آن در مرحله برنامه‌ریزی، اقدام به تصمیم‌گیری می‌نماییم. لذا تعیین می‌کنیم که برای بازه $\Delta(t+1)$ ، تغییرات تعداد ماشین‌های مجازی چگونه خواهد بود. این تصمیم‌گیری بر اساس یک معیار^{۳۸} خاص در مرحله برنامه‌ریزی صورت می‌گیرد. برای مثال می‌توان بهره‌وری پردازنده و یا تعداد درخواست‌ها را ملاک قرار داد. در اینجا بهره‌وری پردازنده را به دلیل دقت بالایی که دارد معیاری برای پیش‌بینی بارکاری در بازه زمانی آتی قرار دادیم.

در ابتدا الگوریتم برای تک تک لایه‌های برنامه کاربردی اجرا می‌شود. سپس در خط ۴ تا ۶ تمامی متغیرهای ارزیابی مقداردهی اولیه می‌شوند. سپس تعداد درخواست‌های وارد شده به لایه ام را با یک سری زمانی جمع می‌کنیم. این سری زمانی برای تخمین تعداد درخواست‌ها در بازه زمانی بعدی به کار می‌رود. همچنین میانه تعداد میلیون دستورالعمل^{۳۹}های هر درخواست را برای لایه ام در مشاهدات n بازه زمانی گذشته تا زمان کنونی را می‌یابیم. در خط ۱۱ با تحلیل‌گر حجم کار از یک مدل رگرسیون خطی (LRM) به منظور پیش‌بینی تعداد درخواست‌ها در بازه زمانی بعدی استفاده می‌کند. رابطه (۱) مدل رگرسیون خطی برای لایه ابر L_i در $t\Delta$ را بیان می‌کند:

$$Y_{t+1} = \beta_1 + \beta_2 \times X_t \quad \text{رابطه (۱)}$$

در رابطه (۲)، t نمونه مشاهده شده است و Yt حجم کار همچون تعداد درخواست‌ها برای لایه ابر L_i است و Xt مقدار زمان واقعی در لحظه‌ای است که هر نمونه گرفته شده است. فرض می‌کنیم n تعداد کل نمونه‌های مشاهده شده است. بنابراین β_1, β_2 با استفاده از حل معادله خطی رگرسیون بر اساس رابطه (۲) تعیین می‌شوند:

38- Metric
39- Million Instructions Per Second (MIPS)

و الگوریتم یادگیری جامع (فراگیر) استوار است، اقدام به تصمیم‌گیری برای مقیاس‌بندی بالا یا پایین بسته به شرایط بهره‌وری پردازنده، در بارکاری آتی می‌نماییم:

$$P_{ScaleUp} = (\text{Predicted Utilization} - 0.8) * 0.5 \quad (۳) \text{ رابطه}$$

$$P_{ScaleDown} = (0.5 - \text{Predicted Utilization}) * 2 \quad (۴) \text{ رابطه}$$

$$P_{Nothing} = 1 - P_{ScaleUp} - P_{ScaleDown} \quad (۵) \text{ رابطه}$$

مطابق رابطه (۳)، چنانچه بهره‌وری بیش از ۸۰ درصد باشد، ScaleUp انجام می‌شود. تخمین بهره‌وری پردازنده در مرحله قبل توسط مدل رگرسیون خطی انجام شد. برای رابطه (۴) نیز فرض کردیم ماشین‌های مجازی که بهره‌وری آن‌ها کمتر از ۵۰ درصد است نیاز به ScaleDown دارند. رابطه (۵) نیز حاصل تفریق احتمالات مقیاس بالا و پایین از احتمال ۱۰۰ درصد می‌باشد.

با استفاده از قضیه بیز احتمال اولیه فرضیه‌ها را به دست می‌آوریم که حاصل آن احتمال یک نوع مقیاس‌بندی در صورت مطلوب بودن برای هر لایه خواهد بود. این احتمال برای هر سه لایه و بار کاری سطح وب، برنامه کاربردی و پایگاه داده محاسبه می‌شود. مفروضات در این بخش بدین شکل می‌باشد:

$$P(\text{Scale Up}|+) = \{0.9, 0.9, 0.87\}$$

$$P(\text{Scale Down}|+) = \{0.97, 0.95, 0.93\}$$

$$P(\text{Nothing}|+) = \{0.8, 0.75, 0.7\}$$

در بیان ساده‌تر در هر بازه زمانی برای هر لایه یکی از سه حالت ScaleUp, ScaleDown و یا Nothing اعمال می‌شود. بنابراین مطابق رابطه (۶) احتمال شرطی برای هر حالت به دست می‌آید:

$$p(h|D) = \frac{p(h)P(D|h)}{P(D)} \quad (۶) \text{ رابطه}$$

$$P(+|\text{Scale Up}) = P(\text{Scale Up}|+)P(+)/P(\text{Scale Up}) \quad (۷) \text{ رابطه}$$

$$P(+|\text{Scale Down}) = P(\text{Scale Down}|+)P(+)/P(\text{Scale Down}) \quad (۸) \text{ رابطه}$$

$$P(+|\text{Nothing}) = P(\text{Nothing}|+)P(+)/P(\text{Nothing}) \quad (۹) \text{ رابطه}$$

در روابط سه‌گانه فوق احتمالات مطلوبیت مقیاس‌بندی‌های مختلف به دست می‌آید. همچنین با توجه

$$\beta_1 = \frac{\sum X_i^2 \sum Y_i - \sum X_i \sum X_i Y_i}{n \sum X_i^2 - (\sum X_i)^2} \quad (۲) \text{ رابطه}$$

$$\beta_2 = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n \sum X_i^2 - (\sum X_i)^2}$$

خط ۱۲ مجموع میانه به دست آمده از تعداد میلیون دستورات عمل‌ها را بر n بازه زمانی تقسیم می‌کند تا میانه برای هر مقطع زمانی به دست آید. بر اساس حجم تخمینی هر درخواست (تعداد میلیون دستورات عمل هر درخواست) و تعداد درخواست‌های تخمینی در بازه زمانی آتی، اقدام به محاسبه حجم کلی تعداد میلیون دستورات عمل‌های دریافتی در بازه زمانی بعدی می‌کنیم (خط ۱۳). در مرحله ۱۴ تخمینی از میزان بهره‌وری در بازه زمانی بعدی را با تقسیم حجم درخواست تخمینی بر ظرفیت کلی منابع لایه نام به دست می‌آوریم. بدین ترتیب میزان بهره‌وری تخمینی همه لایه‌ها را محاسبه و به عنوان خروجی الگوریتم مرحله تحلیل برمی‌گردانیم.

Algorithm3: Pseudocode for Analyze phase

```

1: Input RequestArrivalHistory
   Output predictionMap
2: for (every cloud layer Li in multi layer cloud application at the interval Δt) do
3: begin
4: request_numi=0
5: requestMipsi=0
6: totalMips=0;
7: for t= current_time-n to current_time do
8: requestsNumi,t= Num_of_request_of_layer_i(t)
9: RequestMedianMipsi +=Median_Requests_MIPS(Li,t)
10: end
11: predictedRequestNumberi=Predict number of requests based on local regression
    according to the request numbers in the last n time intervals
12: RequestMedianMipsi= RequestMedianMipsi/n
13: predictedMipsUsagei= predictedRequestNumberi* RequestMedianMipsi
14: predictedUtilizationi=
    predictedMipsUsagei/(Number_Of_VMi*VM_MIPSi*Time_Interval)
15: predictionMap.add(predictedUtilizationi)
16: End
    return predictionMap

```

۳-۲-۴ مرحله برنامه‌ریزی

در اینجا، یک روش مبتنی بر یادگیری بیزی را برای به دست آوردن اقدام بهینه پیاده‌سازی می‌کنیم. ابتدا احتمالات افزایش، کاهش و یا عدم تغییر مقیاس‌بندی را تعیین می‌کنیم. در صورتی که نیاز به منبع بیشتری داشته باشیم مقیاس بالا انجام می‌گیرد. زمانی نیاز به منابع بیشتر احساس می‌شود که میانگین بهره‌وری ماشین‌های مجازی فعلی از حد معینی بالاتر باشد. حال با استفاده از روش پیشنهادی مبتنی بر یادگیری ماشین که بر اصل تئوری بیز

به روابط بالا، احتمال عدم مطلوبیت مقیاس بندی نیز به صورت ذیل به دست می آید:

$$P(-|Scale Up)=1-P(Scale Up|+)P(+)/P(Scale Up) \quad (10)$$

$$P(-|Scale Down)=1-P(Scale Down|+)P(+)/P(Scale Down) \quad (11)$$

$$P(-|Nothing)=1-P(Nothing|+)P(+)/P(Nothing) \quad (12)$$

در انتها مطابق الگوریتم ۴ و روش یادگیری جامع (فراگیر)، ماکزیمم احتمالات محاسبه می شود.

Algorithm4: Pseudocode for planning phase	
1:	Input predictionMap Output Actions
2:	for (every cloud layer L_i in multi layer cloud application at the interval Δt) do
3:	begin
4:	$P(Scale Up)=(\text{predictedUtilization}_i - 0.8)*5$
5:	$P(Scale Down)=(0.5 - \text{predictedUtilization}_i)*2$
6:	$P(Nothing)=1 - P(Scale Up) - P(Scale Down)$
7:	$P(+ Scale Up)=P(Scale Up +)P(+)/P(Scale Up)$
8:	$P(+ Scale Down)=P(Scale Down +)P(+)/P(Scale Down)$
9:	$P(+ Nothing)=P(Nothing +)P(+)/P(Nothing)$
10:	if $P(+ Scale Up) > P(+ Scale Down)$ and $P(+ Scale Up) > P(+ Nothing)$
11:	Scale up a VM in the Layer L_i
12:	else if $P(+ Scale Down) > P(+ Scale Up)$ and $P(+ Scale Down) > P(+ Nothing)$
13:	Scale Down a VM in the Layer L_i
14:	else
15:	Do nothing!

۴-۲-۴- مرحله اجرا

در این مرحله، اقدامات یعنی مقیاس بالا یا پایین که در مرحله برنامه ریزی تعیین شده اند، توسط مدیر VM به منظور درک موثر تامین منبع اجرا می شود و تعامل بین SLA و حداقل هزینه به دست می آید. در الگوریتم ۵ نشان داده شده است که متعادل کننده بار درخواست ها را برای لایه L_i در بازه زمانی به های مناسب بر اساس سیاست های متعادل سازی بار توزیع می کند.

Algorithm5: Pseudocode for Execute phase	
1:	Begin
2:	Execute adaption actions (); /*scale up / scale down by VM manager*/
3:	Dispatch requests to VMs(); /* by Load balancer*/
4:	Sort requests based on deadline
5:	Allocation of requests on VMs of layer I based on BFD algorithm
6:	end

در این مرحله، اعمال انطباقی (Scale Up/Down)

حاصل از مرحله برنامه ریزی، برای تحقق کارایی تامین منابع و همچنین رسیدن به توافق سطح سرویس و کاهش هزینه، توسط مدیر ماشین مجازی اجرا می شوند. درخواست های وارد شده در یک بازه زمانی، به شکل حریصانه^{۴۱} اجرا می شوند. در اینجا تعیین می کنیم که هر

41- Greedy

درخواست مشخص روی کدام ماشین مجازی قرار گیرد. لذا ترتیب درخواست ها از اهمیت بالایی برخوردار است. بنابراین درخواست ها برحسب «ضرب الاجل زمانی»^{۴۲} و به صورت صعودی مرتب می شوند تا از تخطی از توافق سطح سرویس پیشگیری به عمل آید. در طرف مقابل ماشین های مجازی به صورت نزولی و بر اساس میزان بهره وری مرتب می شوند. در نتیجه ظرفیت ماشین های مجازی با ظرفیت کمتر، تکمیل شده و ماشین های مجازی با کمترین میزان بهره وری حذف می شوند. بنابراین تعداد ماشین های مجازی کاهش می یابد.

۵- شبیه سازی روش پیشنهادی و ارزیابی کارایی

در این بخش نتایج حاصل از پیاده سازی روش پیشنهادی برای تامین منابع با استفاده از ترکیب مدل رگرسیون خطی و نظریه بیز را بر اساس سه معیار هزینه، زمان پاسخگویی و میزان بهره وری مورد ارزیابی قرار گرفته است. برای شبیه سازی از ابزار CloudSim استفاده شده و روش پیشنهادی را تحت بارکاری واقعی DPM-RA و Stat-RA FIFA World Cup [۲۱] با روش های مختلف مورد ارزیابی قرار گرفته است. در ادامه ارزیابی کارایی را در سه بخش، تنظیمات شبیه سازی، معیارهای کارایی و ارزیابی نتایج با جزئیات بیشتر ارائه شده است.

۵-۱- تنظیمات شبیه سازی

برای شبیه سازی دقیق و قابل توسعه ابر از ابزار CloudSim استفاده کرده ایم. بدین دلیل که ابزار NetBeans جاوا نوشته شده است، محیط توسعه NetBeans به عنوان محیط برنامه نویسی جاوا انتخاب شده است.

به منظور شبیه سازی راهکار پیشنهادی ما یک برنامه ۳ لایه ابری را در نظر گرفته ایم که هر لایه دارای ویژگی های متفاوتی نسبت به سایر لایه ها است. به عنوان مثال هر لایه نوع ماشین مجازی مجزا و نرخ ورود درخواست های متفاوت را داراست. جدول (۱) مشخصات کلی لایه های

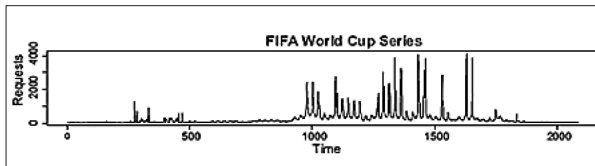
42- Deadline

جدول ۱: مشخصات لایه‌های مختلف برنامه ابری

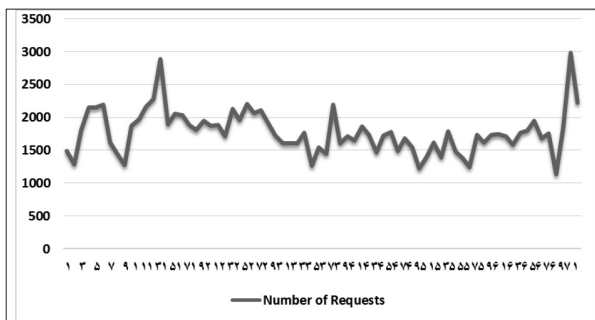
Layer	Avg. Request Arrival Rate	VM Type	Number of Startup VMs
Web	100%	Large	5
Application	70%	Medium	5
Database	40%	Small	5

جدول ۲: مشخصات ماشین‌های مجازی مختلف به کار گرفته شده

VM Type	Large	Medium	Small
CPU (MIPS)	3500	3000	2500
RAM (GB)	4	2	1
Price runtime(\$/hour)	1.28	0.64	0.32
VM Initialization(\$)	1.2	0.6	0.3



نمودار ۱: تعداد تقاضاهای بارهای کاری در محور زمان



نمودار ۲: بار کاری FIFA استفاده شده

مختلف معماری سه لایه‌ای ابری را نشان می‌دهد.

نرخ ورود درخواست‌ها در لایه‌های مختلف متفاوت است. به‌طور دقیق‌تر تعداد درخواست‌های وارد شده در لایه اول مطابق آنچه از بارکاری (FIFA) خوانده می‌شود، می‌باشد اما تعداد درخواست‌های لایه دوم و سوم به ترتیب به‌طور میانگین برابر با ۷۰ و ۴۰ درصد درخواست‌های لایه اول می‌باشد.

همان‌گونه که دیده می‌شود، ماشین مجازی بزرگ به لایه اول (وب) که بیشترین حجم کار را داراست تخصیص داده شده است در حالی که ماشین مجازی با اندازه متوسط به لایه دوم و در نهایت کوچک‌ترین و کم‌قدرت‌ترین نوع ماشین مجازی به لایه آخر (پایگاه داده) که کمترین حجم درخواست‌ها را دریافت خواهد کرد اختصاص یافته است. هر چند تعداد ماشین‌های مجازی اولیه هر سه لایه یکسان (برابر با ۵) است اما به دلیل تفاوت ظرفیت، قدرت و ماهیت ماشین‌های مجازی به‌کار گرفته شده در لایه‌های مختلف، قدرت محاسباتی لایه‌ها از همان ابتدا متفاوت خواهد بود. مشخصات ماشین‌های مجازی مختلف استفاده شده در این سه لایه را در جدول (۲) به‌طور کامل بیان کرده‌ایم.

در این پژوهش به منظور ارزیابی کارایی راهکار پیشنهادی بار FIFA استفاده شده است. در واقع، این بار کاری، حجم درخواست‌های کاربران اینترنتی در محور زمان به این وبگاه مشهور را نشان می‌دهد. داده‌های FIFA شامل ۷ روز کاری و ۱۲۷۸۵۶ درخواست است. شبیه‌سازی در ۲۴ ساعت از ساعات پر درخواست است. نمودار (۱) نمایی از ساختار مجموعه داده واقعی را نشان می‌دهد. ما در این شبیه‌سازی به ترتیب ۶ ساعت از بارکاری به‌صورت ۷۲ بازه زمانی ۵ دقیقه‌ای استفاده کرده‌ایم.

۵-۲ معیارهای کارایی

در این بخش به بررسی کارایی راهکارهای مختلف بر اساس معیارهای مختلف در لایه وب خواهیم پرداخت. تعداد ماشین‌های مجازی، بهره‌وری، زمان پاسخ، هزینه‌ها و میزان سود معیارهای اساسی مورد بررسی ما در این بخش می‌باشند. بدین منظور ما نتایج و تحلیل ارزیابی معیارهای مختلف را برای بار کاری FIFA ارائه داده‌ایم.

ما معیارهای زیر را برای مقایسه رویکرد پیشنهادی با سایر روش‌ها به‌کار گرفته و بدین صورت تعریف می‌کنیم:

الف- بهره‌وری^{۴۲}: بهره‌وری پردازنده برای ماشین‌های مجازی تخصیص داده شده به سرویس ابری S_i در بازه زمانی Δt ، بر اساس تقسیم میزان MIPS تخصیص داده شده ماشین‌های مجازی برای اجرای درخواست‌های سرویس ابری S_i بر کل میزان MIPS ای که به‌طور بالقوه توسط ماشین‌های مجازی برای اجرای درخواست‌های

سرویس ابری در بازه زمانی قابل ارائه است، تعریف میشود و بر اساس رابطه (۱۳) محاسبه می‌شود:

$$U_{S_i}(\Delta t) = \frac{Allocated\ MIPS(\Delta t)}{Total\ MIPS(\Delta t)} \quad (13)$$

ب- ماشین مجازی تخصیص داده شده^{۴۴}: این معیار به صورت تعداد ماشین‌های مجازی تخصیص داده شده به سرویس ابری در بازه زمانی تعریف می‌شود.

پ- تخطی از SLA^{۴۵}: تخطی از SLA برای سرویس ابری S_i در بازه زمانی Δt ، به صورت اختلاف بین زمان پاسخ تعیین شده در SLA کاربر و زمان پاسخ واقعی (تجربه شده) برای درخواست C سرویس ابری S_i در بازه زمانی Δt تعریف می‌شود و بر اساس رابطه (۱۳) و رابطه (۱۴) محاسبه می‌شود:

$$SLAViolation_{S_i}(\Delta t) = \frac{1}{C} \sum_{c=1}^C SLAViolation(request_c) \quad (13)$$

$$SLAViolation(request_c) = \begin{cases} \frac{FT_c - DL_c}{DL_c - AT_c} & FT_c > DL_c \\ 0 & otherwise \end{cases} \quad (14)$$

که AT_c زمان ورود درخواست C کاربر به سیستم، FT_c زمان اتمام درخواست C کاربر و DL_c مهلت اتمام برای درخواست C کاربر می‌باشد.

ت- هزینه کل^{۴۶}: این معیار به عنوان هزینه کل تحمیل شده به تامین کننده SaaS برای پردازش همه درخواست‌های سرویس ابری در بازه زمانی تعریف می‌شود و شامل هزینه ماشین مجازی و هزینه جریمه است که بر اساس رابطه (۱۵) محاسبه می‌شود:

$$Total\ Cost_{S_i}(\Delta t) = VM\ Cost_{S_i}(\Delta t) + Penalty\ Cost_{S_i}(\Delta t) \quad (15)$$

ث- سود^{۴۷}: این معیار به عنوان سود کسب شده توسط تامین کننده SaaS برای پردازش همه درخواست‌های سرویس ابری در بازه زمانی تعریف شده که بر اساس رابطه (۱۶) محاسبه می‌شود:

$$Profit_{S_i}(\Delta t) = Budget_{S_i}(\Delta t) - Total\ Cost_{S_i}(\Delta t) \quad (16)$$

$$Budget_{S_i}(\Delta t) = \sum_{k=1}^K Budget(request_k) \quad (17)$$

که k معرف تعداد درخواست‌های به اتمام رسیده

44- VM allocated
45- SLA violation
46- Total Cost
47- Profit

سرویس ابری S_i در بازه زمانی Δt می‌باشد. بنابراین، بودجه کل سرویس ابری S_i در بازه زمانی برابر مجموع بودجه همه درخواست‌های به اتمام رسیده برای سرویس ابری S_i در بازه زمانی Δt می‌باشد.

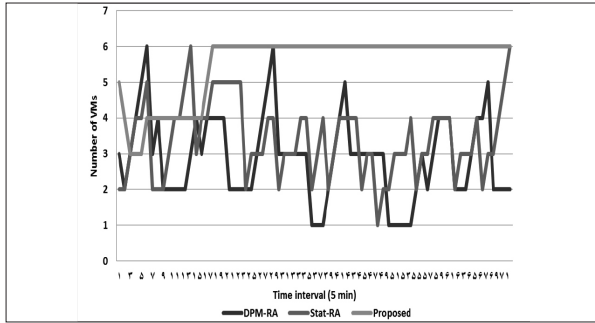
در ابتدا مشاهده می‌شود که درصد میانگین بهره‌وری پردازنده در روش پیشنهادی و در بازه زمانی ۱ تا ۴ بیش از صد درصد می‌باشد. این یعنی تعداد درخواست‌های رسیده به لایه وب در روش ما بیش از ظرفیت منابع موجود می‌باشد. همین مسئله باعث افزایش زمان پاسخ نمودار (۵) به این ترتیب که بازه‌های زمانی کوچکتری در اختیار برنامه‌ها برای اجرا قرار می‌گیرد. از آنجا که بالا رفتن زمان پاسخ موجب پرداخت جریمه می‌شود، لذا به هزینه کلی نیز افزوده می‌شود نمودار (۶). هزینه کلی ما ترکیبی از هزینه اجرای ماشین‌های مجازی، راه‌اندازی آن‌ها و هزینه جریمه می‌باشد. با نگاهی به نمودار (۷) می‌توان دریافت که قسمتی از هزینه کلی مربوط به راه‌اندازی ماشین‌های مجازی می‌باشد و از آنجا که تعداد ماشین‌های مجازی ما نمودار (۳) تا بازه زمانی ۱۶ به طور میانگین کاهش یافته، میزان سود افت چشمگیری را نسبت به روش‌های دیگر نشان نمی‌دهد.

نکته حائز اهمیت در روش پیشنهادی این است که در ابتدای کار به جای افزایش تعداد ماشین‌های مجازی با کاهش روبرو می‌شویم. این یک اشتباه است! از آنجا که روش پیشنهادی بر مبنای احتمالات عمل می‌کند و ما تخمینی از گذشته نداریم، به همین دلیل معمولاً در لایه‌های مختلف و در ابتدای کار ممکن است با تخمین غلط روبرو شویم.

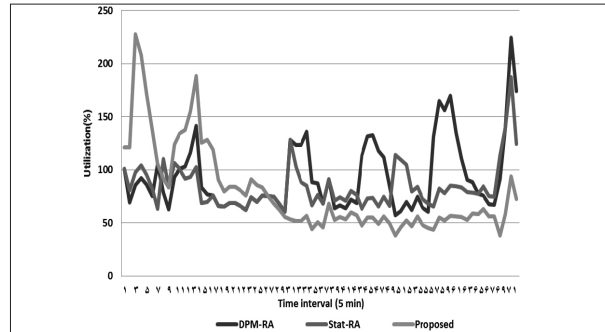
۵-۳ ارزیابی نتایج

۵-۳-۱ ارزیابی معیارهای لایه برنامه کاربردی

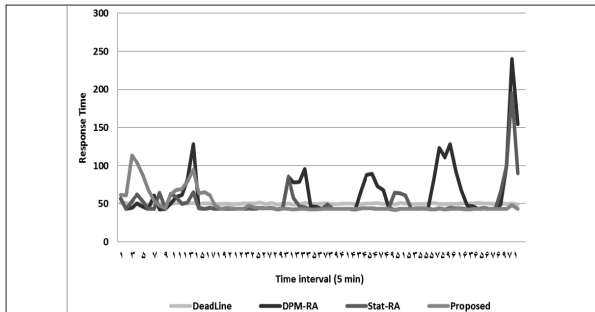
در این بخش به بررسی کارایی روش‌های مختلف بر اساس معیارهای مختلف در لایه برنامه کاربردی خواهیم پرداخت. تعداد ماشین‌های مجازی، بهره‌وری، زمان پاسخ، هزینه‌ها و میزان سود معیارهای اساسی مورد بررسی



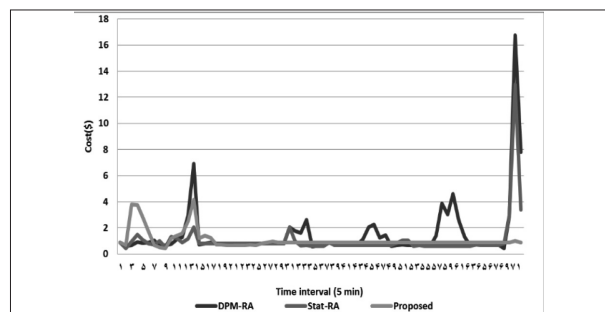
نمودار ۴: میزان بهره‌وری لایه وب



نمودار ۳: تعداد ماشین‌های مجازی لایه وب



نمودار ۶: میزان هزینه ماشین مجازی لایه وب

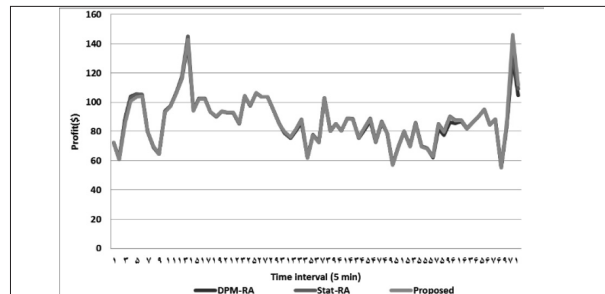


نمودار ۵: میزان زمان پاسخ لایه وب

درخواست‌ها را در بازه زمانی ذکر شده طبق نمودار (۱۰) افزایش می‌دهد که به دلیل تخطی از توافق سطح سرویس منجر به پرداخت جریمه می‌گردد. باید توجه داشت که تعداد ماشین‌های مجازی زیاد همواره نامطلوب نیست. چراکه این اتفاق موجب کاهش بهره‌وری و در نتیجه زمان پاسخ به درخواست می‌شود. کاهش زمان پاسخ موجب کاهش هزینه کلی می‌گردد طبق نمودار (۱۱). که در نهایت بهبود سود را حاصل می‌شود طبق نمودار (۱۲).

۵-۳-۲ ارزیابی معیارهای پایگاه داده

در این بخش اقدام به ارزیابی معیارهای راهکار ارائه شده برای لایه پایگاه داده تحت اجرای بارکاری FIFA خواهیم کرد پیشتر نیز گفته شد که تعداد درخواست‌های لایه پایگاه داده نه تنها کاهش یافته و بلکه می‌تواند متغیرتر از لایه‌های بالاتر باشد (انحراف معیار بیشتر). این قاعده منجر به عملکردی نسبی متفاوت راهکارها نسبت به دو لایه بالایی خواهد شد. نمودار (۱۳) نمودار تعداد ماشین‌های مجازی تخصیص یافته شده راهکارهای مختلف در لایه پایگاه داده را نشان می‌دهد. همان‌گونه که مشخص شده است، به دلیل کاهش تقاضاها در این لایه، به‌طور کلی

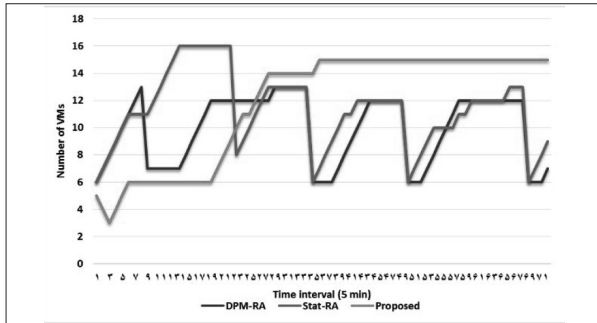


نمودار ۷: میزان سود لایه وب

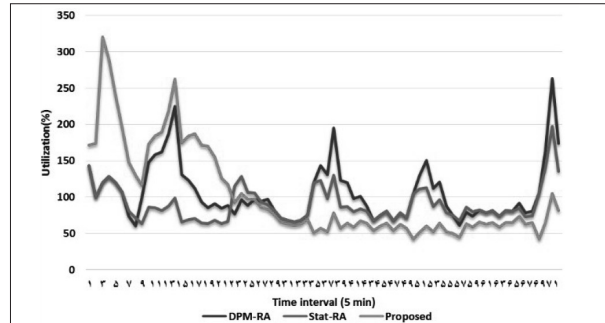
ما در این بخش می‌باشند. بدین منظور ما نتایج و تحلیل ارزیابی معیارهای مختلف را برای بار کاری FIFA ارائه داده‌ایم.

برای بار کاری FIFA تعداد ماشین‌های مجازی لایه دوم نسبت به لایه وب کمی بیشتر است. این اتفاق به دلیل افزایش تقاضاهای دریافتی لایه برنامه کاربردی نسبت به لایه وب افتاده است. به علاوه، سیاست تخصیص تعداد ماشین‌های مجازی مورد استفاده در لایه برنامه کاربردی نیز همانند آنچه در لایه وب اتفاق افتاده تبعیت می‌کند.

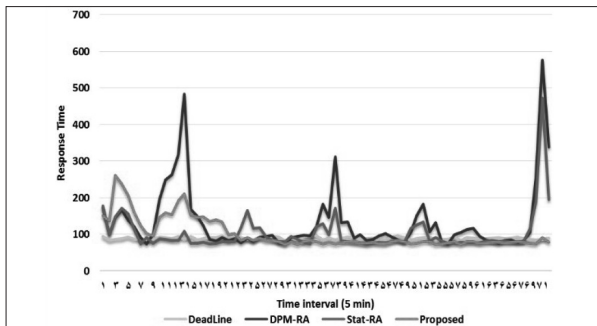
کاهش تعداد ماشین‌های مجازی نسبت به دوراهکار دیگر در بازه زمانی ۱ الی ۲۶، موجب افزایش سهمگین بهره‌وری در نمودار (۹) شده است. همین امر زمان پاسخگویی به



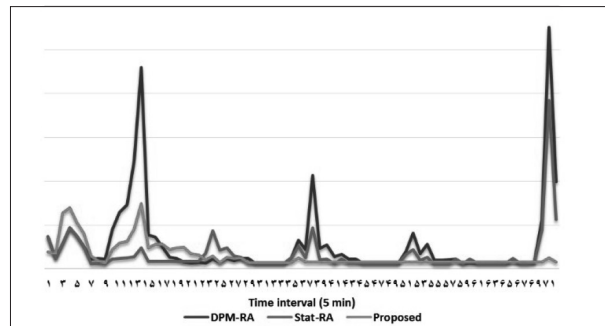
نمودار ۹: میزان بهره‌وری لایه برنامه کاربردی



نمودار ۸: تعداد ماشین‌های مجازی لایه برنامه کاربردی



نمودار ۱۱: میزان هزینه ماشین مجازی لایه برنامه کاربردی

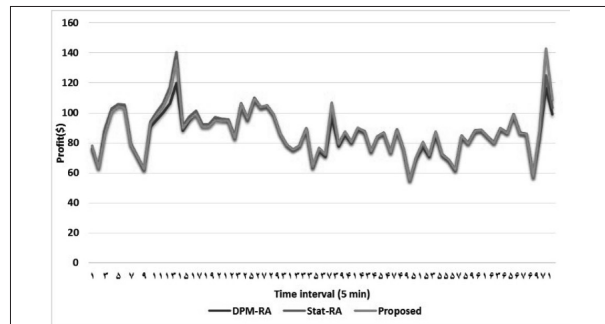


نمودار ۱۰: میزان زمان پاسخ لایه برنامه کاربردی

است. لذا، بهره‌وری روش ماشین یادگیر برای لایه پایگاه داده نیز همانند هر دو لایه بالاتر خود به شکل بهینه‌ای حول حالت تمام بهره‌ور در هر سه بار کاری به شکلی کاملاً بهینه قرار دارد.

نمودار (۱۵) میزان زمان پاسخ درخواست‌ها در لایه پایگاه داده برای بار کاری را نشان می‌دهد. همان‌گونه که مشخص شده است، کاهش نا همگون تقاضاها در این لایه، زمان پاسخ راهکار پیشنهادی را با چالش جدی مواجه نکرده است. از آنجا که تعداد درخواست‌های دریافتی در بازه‌های زمانی مختلف در این لایه تابعی از تعداد درخواست‌های دریافتی در لایه‌های بالاتر است، روند ورود درخواست‌ها در این لایه کمی نامتوازن‌تر از لایه‌های قبلی است. هرچند، زمان پاسخ روش ماشین یادگیر برای لایه پایگاه داده نیز همانند هر دو لایه بالاتر خود به ماشین‌های مجازی سعی در رعایت مهلت زمانی درخواست‌های مشتریان را دارد که نه تنها باعث کاهش هزینه ماشین‌های مجازی و بلکه باعث عدم جلوگیری از افزایش هزینه جریمه می‌شود.

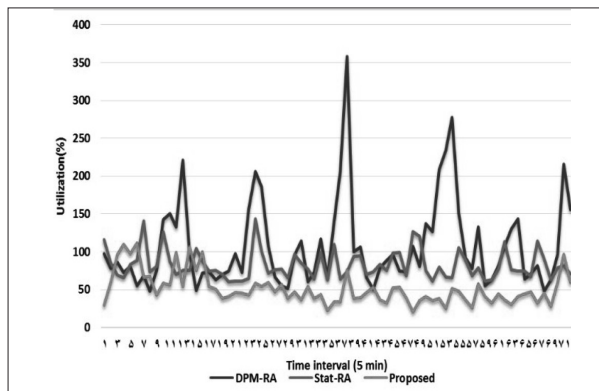
نمودار (۱۶) میزان هزینه در لایه پایگاه داده برای بار کاری FIFA را نشان می‌دهد. راهکار ماشین یادگیر سعی



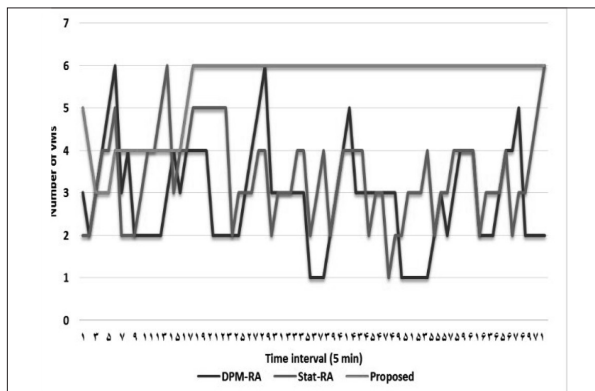
نمودار ۱۲: میزان سود لایه برنامه کاربردی

کمترین تعداد ماشین‌های مجازی مورد استفاده قرار گرفته است. از آنجایی که تعداد درخواست‌های دریافتی در بازه‌های زمانی مختلف در این لایه تابعی از تعداد درخواست‌های دریافتی در لایه‌های بالاتر است، روند ورود درخواست‌ها در این لایه کمی ناهمگون‌تر از لایه‌های قبلی است. با این حال ماشین یادگیر به خوبی توانسته میزان ScaleUp/Down ها را در سطح معینی ثابت نگه دارد.

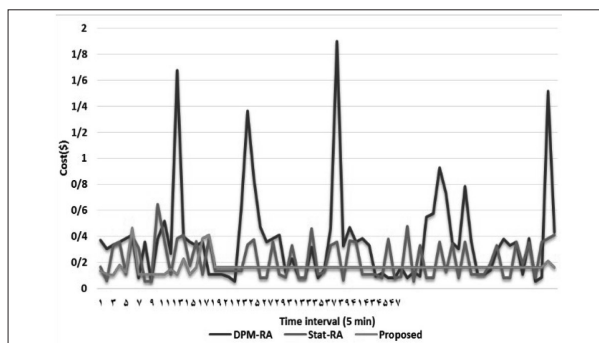
نمودار (۱۴) بهره‌وری ماشین‌های مجازی در لایه پایگاه داده را نشان می‌دهد. همان‌گونه که مشخص شده است، کاهش نامتوازن تقاضاها در این لایه، راهکار پیشنهادی را برخلاف سایر روش‌ها با چالش جدی روبرو نساخته



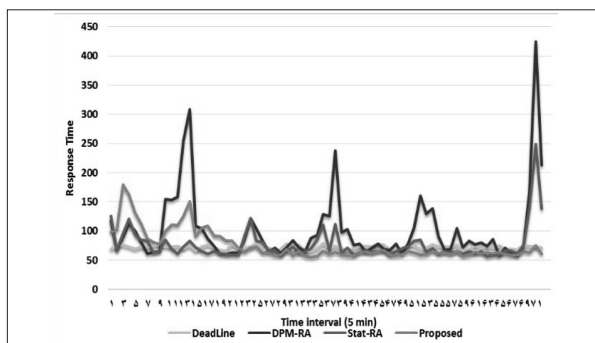
نمودار ۱۴: میزان بهره‌وری لایه پایگاه داده



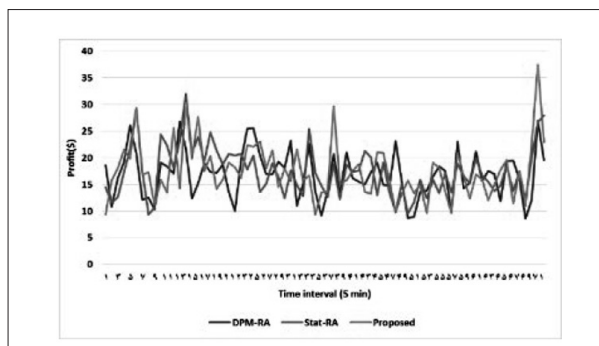
نمودار ۱۳: تعداد ماشین‌های مجازی لایه پایگاه داده



نمودار ۱۶: میزان هزینه ماشین مجازی لایه پایگاه داده



نمودار ۱۵: میزان زمان پاسخ لایه پایگاه داده



نمودار ۱۷: میزان سود لایه پایگاه داده

حاصل تاثیر مثبتی گذاشته است.

ارزیابی میانگین نتایج معیارهای کارایی هر سه لایه

همان‌گونه که در جدول (۴) نشان داده شده است، میانگین تعداد ماشین‌های مجازی مورد استفاده راهکار ارائه شده در بار کاری مورد استفاده بهبود (کاهش) یافته است.

در نتیجه استفاده از تعداد ماشین‌های مجازی بیشتر، بهره‌وری ماشین‌های مجازی راهکار ارائه شده نسبت به دیگر روش‌ها کاهش یافته است. طبق جدول (۵).

از آنجا که تعداد ماشین‌های مجازی استفاده شده جهت

بر کاهش و تثبیت هزینه کلی دارد که نمودار مربوطه بیانگر موفقیت راهکار پیشنهادی نسبت به روش‌های دیگر می‌باشد. نمودار (۱۷) میزان سود در لایه پایگاه داده را نشان می‌دهد. همان‌گونه که مشخص شده است. میزان سود روش ماشین یادگیر برای لایه پایگاه داده بهبود محسوسی (به دلیل انحراف معیار درخواست‌هایش) پیدا کرده است.

ارزیابی شاخص‌های هزینه و سود لایه‌های مختلف

در این بخش به بررسی هزینه‌های راه‌اندازی و اجرای ماشین‌های مجازی، جریمه و میزان سود کلی خواهیم پرداخت. مطابق جدول (۳)، هزینه اجرای ماشین مجازی در روش پیشنهادی بیش از روش‌های دیگر می‌باشد که علت آن افزایش ماشین‌های مجازی توسط ماشین یادگیر می‌باشد. با این حال، هزینه راه‌اندازی راهکار پیشنهادی، به طرز قابل توجهی کمتر از دو روش دیگر شده است که علت آن را باید در کاهش چشمگیر مقیاس‌بندی بالا و پایین جست. همچنین مقدار جریمه نسبت به دو روش دیگر دارای بهبود می‌باشد که همین مهم روی هزینه کل و سود

جدول ۳: هزینه و سود حاصل شده برای راهکارهای مختلف در همه لایه‌ها

روش	هزینه زمان اجرا VM	هزینه راه‌اندازی VM	هزینه کل VM	نرخ جریمه	هزینه جریمه	بودجه درخواستی	هزینه کل	سود کل
MAPE-ML	۱۰۹,۳۲	۷,۴۲	۱۱۶,۷۴	۰,۱۶۷	۸۳,۱۸	۱۴۰۶۹,۰۳	۱۹۹,۹۳	۱۳۸۶۹
DPM-RA	۸۶,۲۶	۲۳	۱۰۹,۲۸	۰,۴۵۷	۲۲۵,۷۹	۱۴۰۳۶,۹۲	۳۳۵,۰۸	۱۳۷۰۱
Stat-RA	۹۵,۰۴	۲۴,۹۷	۱۲۰,۰۲	۰,۱۷۸	۸۸,۴۹	۱۴۰۷۰,۶۴	۲۰۸,۵۲	۱۳۸۶۲

جدول ۵: میانگین بهره‌وری ماشین‌های مجازی مورد استفاده راهکارهای مختلف

روش	MAPE-ML	DPM-RA	Stat-RA
FIFA	۷۷,۴۳	۱۰۲,۳	۸۵,۶۸

جدول ۶: میانگین زمان پاسخ درخواست در راهکارهای مختلف

روش	MAPE-ML	DPM-RA	Stat-RA
FIFA	۱۷۶,۹	۲۲۴,۵	۲۰۹,۵

نوین تحت بار کاری واقعی FIFA و به ارزیابی در عواملی شامل، تعداد ماشین‌های مجازی، بهره‌وری، هزینه مجموع، زمان پاسخ و سود در لایه‌های مختلف پرداختیم و در نهایت روش پیشنهادی MAPE-ML با روش‌های DPM و Stat-RA مقایسه شد.

میانگین تعداد ماشین‌های مجازی مورد استفاده راهکار ارائه شده در بار کاری به‌صورت میانگین، روند افزایشی داشته که در نتیجه آن بهره‌وری پردازنده ماشین‌های مجازی و به تبع آن زمان پاسخگویی به درخواست کاربران کاهش چشمگیری یافته است. هر چند با افزایش تعداد ماشین‌های مجازی برای لایه‌های مختلف در راهکار پیشنهادی، به نظر می‌رسد که با افزایش هزینه روبرو می‌شویم، اما بهبود قطعی زمان پاسخ و پیشگیری از جریمه‌های گزاف، ماشین‌های یادگیر پیشنهادی را با هزینه کلی قابل قبول نسبت به راهکارهای دیگر روبرو می‌سازد. ماشین‌های یادگیر با بهبود پلکانی عملکرد و سعی در ثابت نگه داشتن تعداد ماشین‌های مجازی، نسبت به افزایش سود حاصل اهتمام می‌ورزد.

همچنین راهکار ارائه شده به‌صورت میانگین در هر بازه زمانی باعث بهبود چشمگیری در هزینه‌ها شده است. از آنچه حاصل گردید میانگین حاصل از نتایج شبیه‌سازی،

جدول ۴: میانگین تعداد ماشین‌های مجازی مورد استفاده راهکارهای مختلف

روش	MAPE-ML	DPM-RA	Stat-RA
FIFA	۹,۱	۶,۹	۷,۷۳

سرویس‌دهی به درخواست‌های مشتریان در راهکار ارائه شده بیش از ماشین‌های مجازی سایر راهکارهاست، منجر به کاهش بهره‌وری و تعداد درخواست‌های تخصیص یافته به هر ماشین مجازی می‌شود. در نتیجه کاهش تعداد درخواست‌های تخصیص یافته به هر ماشین مجازی، میانگین زمان پاسخ‌گویی به درخواست‌ها نیز نسبت به سایر راهکارها کاهش داشته است (جدول ۶). بنابراین راهکار پیشنهادی با در پیش گرفتن سیاست کاهش زمان پاسخ، به بهبود سود می‌انجامد.

بر اساس جدول (۷)، اگرچه به نظر می‌رسد با افزایش تعداد ماشین‌های مجازی راهکار ارائه شده متحمل هزینه بیشتری می‌شود اما با پاسخی قابل قبول با میانگین زمان پاسخی مطلوب در مقایسه با سایر روش‌ها و جلوگیری از جریمه‌های غیر معقول، منجر به کاهش هزینه کلی فراهم کننده ابر می‌شود. در نتیجه، راهکار ارائه شده در هر بازه زمانی باعث بهبود قابل توجهی در هزینه‌ها شده است.

در نتیجه آنچه برای سایر معیارها بیان شد، به‌ویژه با کاهش محسوس هزینه راهکار ارائه شده، میانگین میزان سود در هر بازه زمانی حاصل راهکار پیشنهادی نسبت به سایر روش‌ها در بارکاری و در نتیجه بهبود یافته است که در جدول (۸) ارائه شده است.

۶- نتایج و پیشنهادها

در این مقاله به ارائه رویکردی نوین برای تامین کارآمد منابع در محیط‌های ابری و همچنین شبیه‌سازی این روش

Information Systems Frontiers, 16(1), pp.7-18.

[6]. Bi, J., Yuan, H., Tie, M. and Tan, W., 2015. "SLA-based optimisation of virtualised resource for multi-tier web applications in cloud data centres". Enterprise Information Systems, 9(7), pp.743-767.

[7]. Wu, H., Zhang, W., Zhang, J., Wei, J. and Huang, T., 2013. "A benefit-aware on-demand provisioning approach for multi-tier applications in cloud computing". Frontiers of Computer Science, 7(4), pp.459-474.

[8]. Urgaonkar, B., Shenoy, P., Chandra, A., Goyal, P. and Wood, T., 2008. "Agile dynamic provisioning of multi-tier internet applications". ACM Transactions on Autonomous and Adaptive Systems (TAAS), 3(1), p.1.

[9]. RahimiZadeh, K., AnaLoui, M., Kabiri, P. and Javadi, B., 2015. "Performance modeling and analysis of virtualized multi-tier applications under dynamic workloads". Journal of Network and Computer Applications, 56, pp.166-187.

[10]. Beltrán, M., 2015. "Automatic provisioning of multi-tier applications in cloud computing environments". The Journal of Supercomputing, 71(6), pp.2221-2250.

[11]. Singh, S. and Chana, I., 2015. "Q-aware: Quality of service based cloud resource provisioning". Computers & Electrical Engineering, 47, pp.138-160.

[۱۲]. کربلائی مهدی سمانه، قبائی آرانی مصطفی، «رویکرد تامین منبع بهبود یافته برای برنامه‌های کاربردی چندلایه با استفاده از سیستم استنتاج عصبی-مرحله‌ی تطبیق پذیر در محیط رایانش ابری»، مجله علوم رایانشی شماره ۱۳، تابستان ۱۳۹۸

[13]. Khorsand, R., Ghobaei-Arani, M., & Ramezanzpour, M. (2018). "FAHP approach for autonomic resource provisioning of multitier applications in cloud computing environments". Software: Practice and Experience, 48(12), 2147-2173.

[14] Aslanpour, M. S., Ghobaei-Arani, M., & Toosi, A. N. (2017). "Auto-scaling web applications in clouds: A cost-aware approach". Journal of Network and Computer Applications, 95, 26-41.

[15] Guo, J., Li, C., Chen, Y., & Luo, Y. (2020). On-demand resource provision based on load estimation and service expenditure in edge cloud environment. Journal of Network and Computer Applications, 151, 102506.

[۱۶]. اسدی جواد، ناظمی اسلام، «بهبود معماری خودالتیامی در معماری سرویس گرا با استفاده از منطق مرحله‌ی» مجله علوم رایانشی شماره ۴، بهار ۱۳۹۶

[17] Li, W. and Tan, X., 2020. "Locally Bayesian learning in networks". Theoretical Economics, 15(1), pp.239-278.

[18] Alpaydin, E., 2020. "Introduction to machine learning". MIT press.

[19] Barry, D.K., 2003. The Savvy Manager's Guide to Web Services and Service-Oriented Architectures.

[20] Liu, X., Heo, J., Sha, L. and Zhu, X., 2006, April. Adaptive control of multi-tiered web applications using queueing predictor. In Network Operations and Management Symposium, 2006. NOMS 2006. 10th IEEE/IFIP (pp. 106-114). IEEE.

[21]. <http://ita.ee.lbl.gov/html/contrib/WorldCup.html>

جدول ۷: میانگین هزینه در هر بازه زمانی در راهکارهای مختلف

روش	MAPE-ML	DPM-RA	Stat-RA
FIFA	۰,۹۴	۱,۵۲	۰,۹۴

جدول ۸: میانگین سود حاصل از استفاده از راهکارهای مختلف

روش	MAPE-ML	DPM-RA	Stat-RA
FIFA	۶۴,۲۰	۶۳,۴۳	۶۴,۱۷

نشان داد که راهکار ارائه شده در مقایسه با بهترین نتایج به دست آمده از راهکارهای پیشین، منجر به افزایش تعداد ماشین‌های مجازی به میزان ۱۰ درصد با بهبود نرخ مقیاس‌بندی، کاهش ۸ درصدی میانگین بهره‌وری، کاهش ۳ درصدی زمان پاسخ‌دهی، در نتیجه کاهش ۵ درصدی هزینه تمام شده و افزایش ۱ درصدی سود حاصل شده است.

با توجه به نتایج به دست آمده از مقاله حاضر و برای ارزیابی بیشتر نتایج و تکمیل و توسعه این تحقیق، پیشنهاد‌های ذیل ارائه می‌گردد: بهینه نمودن بیشتر روش تامین خودکار منابع در محیط رایانش ابر و اصلاح پارامترها، تعمیم این تکنیک در مجموعه‌های رایانشی بزرگتر و پیچیده‌تر، ارزیابی و بهینه سازی معیارهای دیگر همچون حافظه و پهنای باند، اصلاح حلقه کنترلی MAPE-K در مرحله تحلیل با استفاده از روش ماشین بردار پشتیبان^۸، ترکیب روش Stat-RA و MAPE-ML، ترکیب روش DVM-RA [۷] و MAPE-ML

مراجع:

[1]. Ashouraie, M. and Jafari Navimipour, N., 2015. "Priority-based task scheduling on heterogeneous resources in the Expert Cloud". Kybernetes, 44(10), pp.1455-1471.

[2]. Kupferman, J., Silverman, J., Jara, P. and Browne, J., 2009. "Scaling into the cloud". CS270-advanced.

[3]. Lorigo-Botran, T., Miguel-Alonso, J. and Lozano, J.A., 2014. "A review of auto-scaling techniques for elastic applications in cloud environments". Journal of Grid Computing, 12(4), pp.559-592.

[4]. Jacob, B., Lanyon-Hogg, R., Nadgir, D.K. and Yassin, A.F., 2004. "A Practical Guide to the IBM Autonomic Computing Toolkit", IBM Redbooks.

[5]. Yang, J., Liu, C., Shang, Y., Cheng, B., Mao, Z., Liu, C., Niu, L. and Chen, J., 2014. "A cost-aware auto-scaling approach using the workload prediction in service clouds".

48- Support Vector Machine