

ارائه یک روش مقیاس‌پذیری خودکار به منظور استفاده بهینه از منابع محاسباتی در محیط‌های رایانش ابری

نرگس سادات حسینی مونس

دانشجوی کارشناسی ارشد دانشکده کامپیوتر - دانشگاه علم و صنعت ایران - تهران - ایران
پست الکترونیکی: n_hoseini@comp.iust.ac.ir

مهرداد آشتیانی*

استادیار دانشکده کامپیوتر - دانشگاه علم و صنعت ایران - تهران - ایران
پست الکترونیکی: m_ashtiani@iust.ac.ir

چکیده

محاسبات ابری بسیاری از ارائه‌دهندگان برنامه‌ها را به سمت استقرار برنامه‌های تحت وب بر روی محیط‌های ابری جذب می‌کند. در مقایسه با روش‌های سنتی، رایانش ابری به کاهش هزینه‌ها کمک می‌کند. کاربران محیط‌های ابری می‌توانند در یک بازه زمانی مشخص، تعدادی ماشین مجازی را راه‌اندازی کنند و فقط به اندازه‌ای که از منابع استفاده می‌کنند هزینه پرداخت کنند، که این امر موجب کاهش هزینه‌های تهیه سخت افزارهای حقیقی و پیچیده می‌شود. از آنجایی که میزان درخواست‌ها از برنامه‌های وب در زمان‌های مختلف متفاوت است، تعیین مقدار مناسب منابع ابری مورد نیاز اغلب دشوار است و یکی از چالش‌های مهم در رایانش ابری خواهد بود. گرچه راه‌حل‌های مختلفی برای مدیریت تامین منابع ارائه شده‌است، برای مدیریت مؤثرتر سیستم‌های ذخیره‌سازی مبتنی بر ابر نیاز به روش‌های جدیدتری است. بر این اساس، این کار پژوهشی یک روش پویای مبتنی بر منطق فازی را برای بهبود عملکرد تامین منابع در سیستم‌های ذخیره‌سازی مبتنی بر ابر ارائه می‌کند. نتایج مقایسه

رهیافت پیشنهادی با روش‌های دیگر نشان دهنده عملکرد مثبت این روش در کاهش زمان پاسخگویی به درخواست کاربران و مقدار CPU مصرفی در مرکز داده است. **واژه‌های کلیدی:** مقیاس‌پذیری خودکار، رایانش ابری، استنباط فازی، مدل حلقه OODA

مقدمه

رایانش ابری یک اصطلاح کلی برای هر نوع خدمتی است که شامل ارائه خدمات میزبانی شده از طریق اینترنت است. به عبارتی دیگر، رایانش ابری، در دسترس بودن منابع سیستم رایانه‌ای، به ویژه ذخیره داده (ذخیره‌سازی ابری) و قدرت محاسباتی، بدون مدیریت مستقیم کاربران تعریف شده است. رایانش ابری با ارائه منابع محاسباتی خود به عنوان خدمت، از درگیر کردن کاربران با موضوعاتی همچون خرید، مالکیت و نگهداری از مراکز داده و کارسازهای فیزیکی جلوگیری می‌کند که این موضوع باعث کاهش هزینه‌های کاربران می‌شود [۲]. خدمات رایانش ابری در حال حاضر طیف گسترده‌ای از خدمات را در بر می‌گیرد. این طیف، شامل خدماتی

* نویسنده مسئول

مانند اصول اولیه ذخیره‌سازی، شبکه و قدرت محاسباتی تا پردازش زبان‌های طبیعی و هوش مصنوعی و به‌طور کلی هر خدمتی که نیاز به منابع سخت‌افزاری و فیزیکی دارد، می‌شود. رایانش ابری به دلایل زیادی از جمله صرفه جویی در هزینه، افزایش بهره‌وری، سرعت و کارایی و عملکرد و امنیت، گزینه‌ای محبوب برای شرکت‌ها و ارائه‌دهندگان برنامه‌ها است.

ایجاد برنامه‌های کاربردی ابری با کیفیت بالا یک چالش مهم تحقیقاتی است. در واقع، تخصیص مناسب منابع در شرایطی که نیاز به اعمال خاصیت کشسانی برای مقیاس‌پذیری پویای منابع محاسباتی با توجه به تغییر تقاضا است، یک مسئله مهم و چالش‌برانگیز خواهد بود. پیاده‌سازی یک سازوکار مناسب تخصیص منابع برای به‌دست آوردن توازن میان هزینه و عملکرد بسیار مهم است. نحوه تخصیص منابع باید به گونه‌ای باشد که در عین تامین منابع مورد نیاز برای اجرای درخواست‌های کاربران، مقدار منابع مصرفی و هزینه‌های موجود در حالت کمینه خود باشند و از هر هزینه‌ی اضافی و غیرضروری خودداری شود. برای رسیدن به این هدف، تشخیص شرایط سیستم بسیار حیاتی است زیرا با تشخیص درست شرایط سیستم انتخاب بهترین تصمیم برای تامین منابع عملی می‌شود. با توجه به امکان وجود نوسان در بار ورودی برنامه‌های کاربردی، مطلوب است که یک سازوکار خودکار برای انطباق میزان منابع اختصاص داده شده بر اساس نیازهای سیستم در هر زمان وجود داشته باشد [۳]. مقیاس‌پذیری خودکار یک تکنیک برای تنظیم پویای منابع اختصاص داده شده به برنامه‌های کاربردی با توجه به حجم کار ورودی است. چالش اصلی در مقیاس‌پذیری خودکار، افزایش و یا کاهش منابع با توجه به نوسانات در حجم کار ورودی است. روش‌های مقیاس‌پذیری خودکار بدون دخالت انسان و به‌صورت خودمختار کار می‌کنند [۴]. برای مدیریت موثر برنامه‌های کاربردی، الگوریتم مقیاس‌پذیری درکنار کمینه کردن هزینه باید به رویدادهای سیستم واکنش مناسب

نشان دهد و منابع مورد نیاز برای انجام وظایف سیستم را تامین کند. در مسئله مقیاس‌پذیری خودکار، تکنیک‌های تصمیم‌گیری برای تخصیص منابع به فرایندهای مختلف، به دو دسته واکنشی^۱ و فعال^۲ تقسیم می‌شوند. شیوه واکنشی به‌طور منظم رویدادهایی همچون میزان کار، میزان مصرف پردازنده و صف را مشاهده می‌کند و عملیات مقیاس‌پذیری را با توجه به این رویدادها انجام می‌دهد. اما در روش فعال سعی می‌شود با استفاده از اطلاعات گذشته ازدحام برنامه پیش‌بینی شود و براساس آن تصمیم‌گیری انجام شود. به دلیل سادگی و عملکرد قابل قبول روش‌های واکنشی، بیشتر سازمان‌ها و شرکت‌ها از این روش استفاده می‌کنند. مشکل رویکرد واکنشی این است که آستانه و تصمیمات مقیاس‌پذیری، با توجه به حجم کار و شرایط سامانه تغییر نمی‌کند و ثابت است که این موضوع می‌تواند مانع انتخاب بهترین تصمیم در تامین منابع شود. یک روش برای حل این مشکل، اضافه کردن پویایی به رویکردهای واکنشی بوده به این صورت که در بازه‌های زمانی مشخص، شرایط سامانه مشخص شود و متناسب با شرایط سامانه در آن بازه زمانی مشخص، رویکرد واکنشی به تصمیم‌گیری بپردازد.

در این کار پژوهشی ما یک رویکرد مقیاس‌پذیری واکنشی را پیشنهاد می‌کنیم که مقیاس‌پذیری را بر اساس مشخصه‌های سامانه به‌صورت پویا انجام می‌دهد. برای این منظور، منطق فازی بسیار کاربردی به نظر می‌رسد زیرا برخلاف روش‌های یادگیری ماشین برای ساختن مدل به داده و زمان آموزش طولانی نیاز ندارد و می‌تواند عملکرد مناسبی داشته باشد. در روش پیشنهاد شده تلاش شده‌است با استفاده از برخی از ویژگی‌های سامانه مانند زمان پاسخگویی به درخواست‌ها و مقدار مصرف پردازنده توسط ماشین‌های مجازی، وضعیت سامانه مشخص شود و با توجه به وضعیت حال حاضر سامانه، متغیرهای مورد نیاز در تصمیم‌گیری برای مقیاس‌پذیری

1- Reactive

2- Proactive

مانند قوانین فازی، به صورت پویا و متناسب با وضعیت حال حاضر استفاده شوند. استفاده از رویکرد واکنشی در کنار پویایی اضافه شده به آن به ما کمک می‌کند که علاوه بر سرعت در تصمیم‌گیری، به دقت مناسب و قابل قبول متناسب با شرایط، دست پیدا کنیم.

در مسئله مورد بررسی، مرکز داده محیط رایانش ابری شامل تعدادی ماشین فیزیکی است. هر ماشین مجازی دارای تعدادی هسته پردازنده بوده و تمامی ماشین‌های مجازی در مرکز داده از نظر قدرت محاسباتی یکسان هستند. برای هر کاربر، تعدادی ماشین مجازی در نظر گرفته شده است که وظیفه اجرای کارهای ارسالی از طرف کاربران را بر عهده دارند. کارهای ورودی از طرف کاربران از یکدیگر مستقل می‌باشند و منابع اختصاص داده شده به آن‌ها در طول اجرا ثابت می‌باشد. هر کار ارسالی از طرف کاربر دارای طول مشخصی می‌باشد که با توجه به قدرت محاسباتی منبع اختصاص داده شده به کار ورودی، می‌توان یک مهلت اجرای تخمینی برای هر وظیفه از طرف کاربر در نظر گرفت. روش پیشنهاد شده، یک روش مقیاس‌پذیری واکنشی فازی است که منطبق با مدل حلقه OODA، با در نظر گرفتن ویژگی‌های سیستم، مانند مدت زمان پاسخگویی به درخواست کاربران و میانگین مقدار پردازنده مصرف شده، بهترین قوانین فازی را از مجموعه قوانینی که از پیش برای حالت‌های مختلف سیستم تعریف شده‌است، انتخاب کرده و عملیات مقیاس‌پذیری را به صورت خودکار برای تامین منابع مورد نیاز کارهای درخواستی توسط کاربران که مستقل و غیرقابل پیش‌بینی هستند، انجام می‌دهد. قوانین تعریف شده می‌توانند به سادگی «اگر مدت زمان پاسخگویی زیاد بود، تعداد منابع محاسباتی افزایش پیدا کند»، باشد. روش پیشنهاد شده بدون نیاز به جمع‌آوری و ذخیره حجم داده زیاد، به دنبال کاهش متغیرهایی همچون مدت زمان پاسخگویی و میانگین قدرت محاسباتی مصرفی منابع است.

روش پیشنهادی از نظر عملکرد در مقایسه با دیگر

روش‌های واکنشی معرفی شده، با توجه به وجود مؤلفه پویایی، عملکرد بهتری دارد و از نظر قدرت محاسباتی و حافظه مورد نیاز برای اجرای این روش، در مقایسه با روش‌های فعال، کم‌هزینه و مناسب‌تر است. نوآوری و دستاوردهای این کار پژوهشی به طور خلاصه شامل موارد زیر است:

۱. خودکارسازی فرآیند مقیاس‌پذیری بدون دخالت انسان و در نتیجه دستیابی به تصمیم‌گیری بهینه برای بهبود مؤلفه‌های ارزیابی
 ۲. عدم نیاز به جمع‌آوری حجم بالای داده تاریخی.
 ۳. عدم نیاز به صرف زمان آموزش.
 ۴. پویایی مضاعف در مراحل مختلف الگوریتم.
- در ادامه پس از مقدمه اولیه گفته شده در این بخش، در بخش دوم به معرفی و توضیح ادبیات و مفاهیم اولیه مورد استفاده در این کار پژوهشی مانند مفاهیم مرتبط با مقیاس‌پذیری خودکار و توضیح عملکرد سیستم‌های استنباط فازی پرداخته می‌شود. در بخش سوم تعدادی از کارهای انجام شده و پژوهش‌های اخیر در حوزه مقیاس‌پذیری خودکار معرفی می‌شوند و نقاط مثبت و منفی هر کدام توضیح داده می‌شود. در بخش چهارم، جزئیات رهیافت پیشنهادی تشریح می‌شود. در بخش پنجم، ارزیابی رهیافت پیشنهادی در مراحل متفاوت، مانند مرحله مقایسه با روش‌های دیگر و یا مرحله تحلیل حساسیت، انجام می‌شود و در بخش آخر، بخش ششم، به جمع‌بندی و نتیجه‌گیری از پژوهش انجام شده و ارائه پیشنهاد برای کارهای آینده، پرداخته می‌شود.

۲- مفاهیم پایه

در این بخش به صورت خلاصه به معرفی برخی از مفاهیم پایه‌ای برای درک بهتر این مقاله پرداخته می‌شود.

۲-۱- مقیاس‌پذیری خودکار

یک سامانه که از ویژگی کشسانی محیط ابری با اختصاص مقدار مناسب منابع به حجم کار برنامه‌ها،

با نظارت دستی اندک و یا بدون نظارت دستی استفاده می‌کند، می‌تواند به‌عنوان یک سامانه مقیاس‌پذیر خودکار در نظر گرفته شود. مقیاس‌پذیری خودکار یک روش برای تنظیم پویای منابع محاسباتی اختصاص داده شده به برنامه‌های کاربردی با توجه به حجم کار ورودی است [۵]. چالش اصلی در مقیاس‌پذیری خودکار افزایش و یا کاهش منابع با توجه به نوسانات در حجم کار ورودی خواهد بود. روش‌های مقیاس‌پذیری خودکار بدون دخالت انسان و به‌صورت خودمختار کار می‌کنند [۴].

۲-۱-۱- انواع مقیاس‌پذیری

عملیات مقیاس‌پذیری به دو دسته واکنشی و فعال تقسیم می‌شوند [۶]. زمان انجام عملیات مقیاس‌پذیری، نشان دهنده روش اتخاذ شده برای مقیاس‌پذیری است. در روش واکنشی، پس از به وجود آمدن تغییرات در بار ورودی و تحلیل وضعیت فعلی سیستم، منابع اختصاص داده می‌شوند در حالی‌که در روش فعال، ابتدا سعی می‌شود تغییرات بار ورودی با استفاده از داده‌هایی که از عملکرد سیستم در زمان‌های گذشته جمع‌آوری شده است، پیش‌بینی شود و براساس این پیش‌بینی اختصاص منابع اتفاق می‌افتد. امروزه در بسیاری از روش‌های جدید پیشنهاد شده، ترکیبی از دو روش گفته شده برای رسیدن به نتایج بهتر استفاده می‌شود.

۲-۱-۲- سیاست‌های مقیاس‌پذیری

مقیاس‌پذیری می‌تواند در دو حالت افقی و یا عمودی اتفاق بیفتد [۷]. مقیاس‌پذیری افقی به معنای افزایش و یا کاهش تعداد منابع موجود در مرکز داده و مقیاس‌پذیری عمودی به معنای افزایش و یا کاهش ظرفیت و قدرت منابع موجود در مرکز داده است. ارائه‌دهندگان خدمات ابری اغلب از مقیاس‌پذیری افقی برای تامین خودکار منابع استفاده می‌کنند به این صورت که یک منبع قابل تنظیم در اختیار کاربر قرار می‌گیرد که در آن کاربر می‌تواند ماشین‌های مجازی را با تعیین حافظه، هسته، پهنای باند و

دیگر ویژگی‌های لازم پیکربندی کند و با توجه به نیاز خود، تعدادی ماشین‌مجازی اضافه و یا کم کند. در مقیاس‌پذیری عمودی، کاربر می‌تواند منابع ماشین‌های مجازی همانند CPU، حافظه و پهنای باند شبکه را که در اختیار دارد، مجدداً پیکربندی کند.

۲-۱-۳- رویکردهای مقیاس‌پذیری

در روش‌های مقیاس‌پذیری خودکار، رویکرد مقیاس نشان‌دهنده متغیر یا متغیرهایی است که تمرکز روش بر روی بهبود آن است. در مقیاس‌پذیری، چشم‌انداز روش می‌تواند بهبود یک یا چند عامل مانند هزینه، بهینه‌سازی منابع و بهبود کیفیت خدمات باشد. مشخص کردن رویکرد مقیاس‌پذیری به افراد کمک می‌کند تا روش مقیاس‌پذیری خودکار را با توجه به نیاز مشخص و خاص خود انتخاب کنند.

۲-۱-۴- معیارهای عملکرد

در روش‌های متفاوتی که وجود دارد، مهم‌ترین مرحله، ارزیابی مدل و روش پیشنهاد شده و بررسی پارامترهایی است که روش در بهبود آن‌ها عملکرد مثبتی داشته است. پارامترها و معیارهای متنوعی مانند زمان پاسخ، میزان استفاده از CPU و میزان درخواست ورودی برای بررسی عملکرد سیستم استفاده می‌شوند. انتخاب روش مناسب مقیاس‌پذیری با توجه به معیارهای عملکرد روش، یک مرحله ضروری برای افزایش قابلیت اطمینان کاربرانی است که برنامه‌های خود را به محیط‌های رایانش ابری منتقل می‌کنند.

۲-۲- حلقه OODA

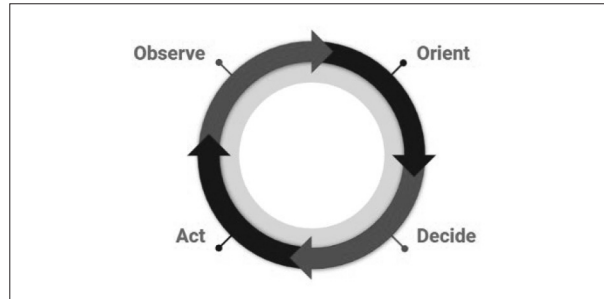
برای مدیریت موثر برنامه‌های کاربردی، الگوریتم مقیاس‌پذیری باید به رویدادهای سیستم واکنش مناسب نشان دهد. با توجه به این‌که فرایند مقیاس‌پذیری خودکار با مدل حلقه OODA (شکل ۱) مطابقت دارد، در بسیاری از الگوریتم‌ها از این مدل برای مقیاس‌پذیری خودکار استفاده می‌شود. این حلقه از چهار مرحله شامل مشاهده^۲،

3- Observation

۳ - مروری بر کارهای مرتبط

تاکنون روش‌های متفاوتی در زمینه مقیاس‌پذیری خودکار پیشنهاد شده است. این روش‌ها در گذشته بیشتر با استفاده از روش‌های واکنشی مانند روش‌های مبتنی بر آستانه و روش‌های تحلیل سری‌های زمانی، پیاده‌سازی می‌شدند اما امروزه بیشتر روش‌ها به تکنیک فعال روی آورده‌اند و اغلب با استفاده از روش‌های یادگیری ماشین، یادگیری تقویتی و روش‌های دیگر در زمینه هوش مصنوعی، پیاده‌سازی می‌شوند. مشکل روش‌های فعال نیاز آن‌ها به حجم زیاد داده، قدرت محاسباتی قابل قبول و زمان طولانی برای اجرا است. البته امروزه با پیشرفت فناوری، استفاده از این روش‌ها بسیار ساده و قابل قبول شده است. از طرفی مشکل اساسی روش‌های واکنشی، ایستا بودن این روش‌ها است. در روش‌های مقیاس‌پذیری خودکار به صورت واکنشی، الگوریتم، همواره اطلاعات محیط رایانش ابری را بررسی می‌کند و پس از بررسی اطلاعات به تصمیم‌گیری درباره تامین منابع می‌پردازد. برای پیاده‌سازی الگوریتم‌های واکنشی مقیاس‌پذیری خودکار از روش‌هایی مانند استنباط فازی [۱]، روش‌های مبتنی بر آستانه [۹] و روش‌های مبتنی بر تطبیق الگو [۱۰]، استفاده می‌شود.

روش‌های مبتنی بر آستانه از مهم‌ترین روش‌های مقیاس‌پذیری خودکار با استفاده از تکنیک واکنشی می‌باشند. در روش پیشنهاد شده توسط آقای Hung و همکارانشان [۹]، یک روش مقیاس‌پذیری خودکار مبتنی بر آستانه پیشنهاد شده است. در روش پیشنهاد شده، درخواست‌های کاربران توسط متعادل‌کننده بار میان ماشین‌های مجازی موجود در هر خوشه، توزیع می‌شوند. هر ماشین مجازی شامل تعدادی نشست^۴ است و وظیفه اجرای کارهای ورودی بر عهده آن‌ها خواهد بود. در این روش، پارامتر مورد بررسی، تعداد نشست‌های فعال در هر ماشین مجازی است. در صورتی که تعداد نشست‌های فعال در یک ماشین مجازی از آستانه بالایی مشخص شده بیشتر



شکل ۱: مدل حلقه OODA [۱]

جهت‌گیری^۵، تصمیم‌گیری^۶ و عملیات^۷ تشکیل شده است.

۲-۳- استنباط فازی

نظریه منطق فازی برای استفاده از دانش متخصص در قالب قوانین فازی برای تعیین یک اقدام مناسب و بهینه طراحی شده است. منطق فازی می‌تواند در هرس فضای حالت مسئله نیز تاثیرگذار باشد. استنباط فازی با مقایسه مقادیر واقعی و مقادیر مرجع عوامل موثر در سیستم، شرایط را کنترل می‌کند. منطق فازی برای استفاده در محیط‌های تصادفی و پویا با مجموعه‌ای از عدم قطعیت‌ها مانند شبکه‌های کامپیوتری، سیستم‌های توزیع شده و محیط‌های ابری مناسب است و برای حل مشکلات تصمیم‌گیری در حوزه مختلف مانند اینترنت اشیا، داده‌کاوی، رایانش ابری، شبکه‌های اجتماعی و پردازش تصویر، مورد استفاده قرار گرفته است. یک سیستم استنباط فازی با فازی‌سازی معیارهای ورودی کار می‌کند. تعداد معیارهای ورودی می‌تواند یک، دو، سه و حتی بیشتر باشد. پس از دریافت اطلاعات ورودی، یک موتور استنباط به تصمیم‌گیری بر اساس داده‌های ورودی و پایگاه قاعده می‌پردازد. سیستم استنباط فازی از قسمت‌های پایگاه قاعده، واحد فازی‌سازی، تابع عضویت، واحد استنباط فازی و واحد خنثی‌سازی فازی تشکیل شده است [۳]. برای انجام استنباط فازی روش‌های متفاوتی از جمله دو روش ممدانی و سوگنو^۷ وجود دارد که در این کار تحقیقاتی از روش ممدانی استفاده شده است [۸].

4- Orientation

5- Decision

6- Action

7- Sugeno

باشد، الگوریتم، دستور افزایش تعداد ماشین‌های مجازی را اعلام می‌کند و این افزایش به اندازه یک ماشین مجازی خواهد بود. در شرایطی که تعداد نشست‌های فعال در یک ماشین مجازی از آستانه پایینی مشخص شده کمتر باشد، آن ماشین مجازی بیکار و بدون استفاده است. برای کاهش تعداد ماشین‌های مجازی در صورتی که بیش از دو ماشین مجازی در یک خوشه بیکار باشند، الگوریتم دستور تخریب و کاهش یکی از ماشین‌های مجازی را می‌دهد. با توجه به سادگی الگوریتم بیان شده، این الگوریتم از نظر اجرا مدت زمان کمی نیاز دارد اما این الگوریتم نیز همانند بسیاری از روش‌های مبتنی بر آستانه دیگر، از مرزهای ایستا استفاده کرده که این موضوع مانع از رسیدن به عملکرد مطلوب می‌شود. زیرا شرایط سیستم همواره ثابت نیست و به صورت پویا در تغییر بوده و با در نظر گرفتن یک آستانه ثابت، مهم‌ترین ویژگی سیستم‌های رایانش ابری که پویایی آن است را در نظر نگرفته‌ایم.

آقای Liao و همکاران ایشان [۱۴] نیز از یک روش مبتنی بر آستانه برای الگوریتم مقیاس‌پذیری خودکار استفاده کردند که تفاوت اصلی این روش با روش قبلی در پویا بودن مرزها و آستانه‌های آن است. هدف اصلی روش ارائه شده، مقیاس‌پذیری خودکار برای خدمات وب آمازون است. در این روش، متغیری که آستانه برای آن مشخص شده، مقدار استفاده از منبع پردازنده مرکزی است. ایده اصلی در این روش، پویایی مرزهای آن است. در توضیح این روش آمده است که زمانی که درصد استفاده از CPU بالا بوده و با بررسی شرایط محیط متوجه افزایش بار کاری هر ماشین مجازی می‌شویم، احتمال این که در تامین منابع دچار کمبود شویم بالا می‌رود و به همین علت و برای جلوگیری از کمبود منابع، مرز بالایی با هر افزایش بار، مقداری کاهش پیدا می‌کند. به صورت مشابه، برای حالتی که با بررسی شرایط هر ماشین مجازی متوجه می‌شویم که بار ورودی در حال کاهش است، احتمال تامین منابع بیش از حد مورد نیاز وجود داشته و در نتیجه، آستانه

بالایی با هر کاهش بار، مقداری افزایش پیدا می‌کند. مقدار تغییر آستانه، چه در کاهش و چه در افزایش به اندازه ۰/۵ درصد مقدار حال حاضر آستانه است.

برای آستانه پایانی نیز تغییرات پویا صادق خواهد بود. تغییرات مرز پایینی مشابه مرز بالایی بوده به این صورت که در حالتی که در مرحله افزایش بار باشیم، مرز پایینی افزایش پیدا می‌کند و در صورتی که در مرحله کاهش بار باشیم، آستانه پایینی کاهش پیدا خواهد کرد. نکته‌ای که در مرز پایینی وجود دارد این است که در شرایطی که بیش از ۵۰ درصد از ماشین‌های مجازی در حال انجام وظایف باشند، برای جلوگیری از نوسان در مقیاس‌پذیری و کاهش و افزایش متناوب منابع، مرز پایینی در کم‌ترین حالت خود قرار می‌گیرد. همان‌طور که مشخص است استفاده از این روش می‌تواند در بهبود عملکرد سیستم بسیار موثر باشد. زیرا با تعیین مرزهای پویا، الگوریتم، هم با ویژگی پویایی سیستم سازگار شده و عملکرد مناسبی دارد و هم از نوسان‌های نامطلوب در مقیاس‌پذیری جلوگیری می‌کند.

روش‌های مقیاس‌پذیری مبتنی بر استنباط فازی، از منطق و سیستم استنباط فازی برای تصمیم‌گیری در مورد نحوه تامین منابع استفاده می‌کنند. در این روش‌ها پس از جمع‌آوری داده مورد نیاز برای استنباط فازی، مقادیر ورودی به یک سیستم استنباط فازی داده می‌شود و خروجی مورد نظر تحویل گرفته می‌شود. در روش‌های متفاوت وابسته به شرایط خاص هر مسئله، خروجی سیستم استنباط فازی متفاوت است. اما به‌طور کلی خروجی، نحوه مقیاس‌پذیری را مشخص می‌کند. مزیت اصلی این روش نسبت به روش آستانه‌گذاری، پیچیدگی بیشتر آن بوده که از پارامترهای بیشتر و پیچیده‌تری برای تصمیم‌گیری استفاده می‌کند.

روش پیشنهاد شده توسط Tseng و همکاران [۱۱] یک روش مبتنی بر استنباط فازی است. در این روش، تعدادی متغیر به‌عنوان ورودی در نظر گرفته شده و به سیستم استنباط فازی داده می‌شود تا خروجی مورد نیاز تولید

شود. برای استفاده از سیستم فازی ابتدا باید اصطلاحات فازی و توابع عضویت متغیرها مشخص شود. در این روش از تابع توزیع گاوسی برای توابع عضویت متغیرها استفاده شده است. نکته قابل توجهی که وجود دارد این است که تابع عضویت برای هر متغیر در شرایط مختلف متفاوت است. به این معنا که برای هر متغیر سه حالت زیاد، متوسط و کم وجود دارد که تابع عضویت هر متغیر متناسب با شرایطی که دارد تعریف و استفاده می‌شود. در این رهیافت، متغیرهای ورودی مورد بررسی مقدار CPU مصرفی، حافظه و شبکه و متغیر نهایی و خروجی، متغیر مقیاس است. برای هر کدام از متغیرهای خروجی و ورودی، سه اصطلاح فازی کم، متوسط و زیاد تعریف شده است. با توجه به این که این روش از منطق فازی برای تصمیم‌گیری استفاده می‌کند، عملکرد قابل قبول و مناسبی دارد اما در نقطه مقابل نقاط ضعفی نیز برای این روش می‌توان در نظر گرفت از جمله این که پایگاه قاعده دچار مشکل ایستایی بوده و یک پایگاه قاعده با یک سری از قوانین در تمامی حالات استفاده می‌شوند که این موضوع می‌تواند در عملکرد روش اثر مخربی داشته باشد.

در کار پژوهشی انجام شده توسط آقای قبائی و همکاران ایشان [8] نیز از سیستم استنباط فازی برای مقیاس‌پذیری استفاده شده است. این روش همانند رهیافت پیشنهادی ارائه شده، از مدل حلقه OODA برای مدیریت بهتر استفاده می‌کند. سیستم مورد بررسی در این روش، دارای تعدادی گره داده بوده که وظیفه اجرای وظایف محول شده از طرف کاربران را بر عهده دارند. در این روش، متغیرهای ورودی به سیستم استنباط فازی، دو متغیر VPT به معنای تعداد بازدید از یک بلوک داده در یک بازه زمانی و AVT به معنای میانگین زمانی که برای بازدید یک بلوک داده صرف می‌شود، هستند. برای فازی‌سازی این دو متغیر، ابتدا دو مرز بالایی و پایینی برای هر کدام در نظر گرفته می‌شود و سپس فازی‌سازی آن‌ها به کمک سه اصطلاح فازی کم، میانه و زیاد انجام می‌شود. نقاط ابتدایی

و انتهایی اصطلاح فازی میانه، دو آستانه بالایی و پایینی هر متغیر خواهند بود.

مرزهایی پایینی و بالایی هر متغیر با تغییر شرایط سیستم، تغییر می‌کند و تغییرات آن وابسته به تعداد دفعاتی است که سیستم دچار نقض توافقنامه در سطح سیستم شده باشد. زمانی یک وظیفه و درخواست دچار نقض توافقنامه در سطح سیستم می‌شود که مدت زمان اجرای آن وظیفه از مدت زمان مشخص شده برای اجرای آن وظیفه توسط کاربر، بیشتر شده باشد. این رهیافت، برای استنباط فازی خود از روش ممدانی استفاده کرده است. در مرحله آخر نیز با مشخص شدن مقدار متغیر خروجی، یک منبع محاسباتی به محیط رایانش اضافه یا کم می‌شود. با بررسی این روش، مشخص می‌شود که این روش نسبت به روش قبل پویایی بیشتری دارد (مانند پویایی مرزها در اصطلاحات زبانی متغیرها و تغییر آن‌ها متناسب با توافقنامه سطح سیستم)، اما همچنان می‌توان به این روش پویایی بیشتری اضافه کرد. به طور مثال، مانند روش قبلی، به جای استفاده از یک پایگاه قاعده ثابت، پویایی به پایگاه قاعده نیز اضافه شود و یا تعداد منابعی که به محیط محاسباتی اضافه و یا کم می‌شوند نیز یک عدد ثابت نباشد و متناسب با شرایط سیستم تغییر کند.

در روش‌های مقیاس‌پذیری مبتنی بر تطبیق الگو، الگوریتم مقیاس‌پذیری در ابتدا بدون هیچ اطلاعات اولیه در مورد سیستم، شروع به کار می‌کند. با گذر زمان، الگوریتم، اطلاعات مورد نیاز در مورد نحوه تغییر بار ورودی و تغییرات متغیرهای موثر در مقیاس‌پذیری را جمع‌آوری می‌کند. پس از جمع‌آوری داده‌های مورد نیاز، الگوریتم به دنبال الگوهای مشترک و تکراری در تغییرات پارامترهای سیستم می‌گردد. برای پیدا کردن الگوهای مناسب و درست، معمولاً در ابتدا یک مجموعه نامزد به عنوان الگوهای تکرار شونده در سیستم، در نظر گرفته می‌شوند. سپس، با معیارهایی که از پیش تعریف شده است، درستی هر کدام از الگوها بررسی می‌شود و

9- Pattern matching

از مجموعه نامزد، وابسته به الگوریتم ارائه شده، یک یا چند الگو به عنوان الگوی تکراری در نظر گرفته می‌شوند. در نهایت، الگوریتم در شرایطی که متوجه وقوع این الگو شود، دستور متناسب با شرایطی که از وقوع این الگو پیش‌بینی می‌کند را صادر می‌کند.

در روش پیشنهاد شده توسط Nowick و دیگران، از روش تطبیق الگو برای مقیاس‌پذیری خودکار استفاده شده است [۱۰]. در الگوریتم پیشنهادی، از دو قسمت برای تخمین شرایط سیستم و در نتیجه تصمیم‌گیری در رابطه با نحوه مقیاس‌پذیری، استفاده شده است. در قسمت اول، الگوریتم برخی از فاکتورهای تاثیرگذار بر بار ورودی را مورد بررسی قرار می‌دهد. از جمله این فاکتورها که در بار ورودی تاثیرگذارند می‌توان به، بازه زمانی روز، تعطیلات و روزهای آخر هفته اشاره کرد. قسمت دوم یک الگوریتم تطبیق الگوی دوفازی^{۱۰} را پیاده‌سازی می‌کند.

در این روش، پس از جمع‌آوری داده، به تحلیل سری‌های زمانی پرداخته می‌شود. برای به‌دست آوردن تناوب در سری زمانی، از خودهمبستگی^{۱۱} استفاده شده است که نشان‌دهنده همبستگی متقاطع یک سیگنال با خودش است. پس از به‌دست آوردن الگوهای متناوب، مرحله بعدی، یافتن الگوهای مشابه با الگوهای موجود در داده‌های تاریخی است. در حقیقت در این قسمت، الگوریتم می‌خواهد تعیین کند که در گذشته چه زمانی نوسانات و مقادیر حجم کار مشابه ثبت شده است. شباهت بر اساس مقدار حجم کار و همچنین نوسانات الگو (اعم از افزایش، کاهش یا پایدار بودن حجم کار) تعریف می‌شود. پس از به‌دست آوردن الگوهای متناوب، الگوریتم با استفاده از روش Relief، به ارزیابی الگوهای (ویژگی‌ها) به‌دست آمده می‌پردازد و سپس از این داده‌ها برای مرحله یادگیری ماشین استفاده می‌کند. در مرحله یادگیری، الگوریتم با استفاده از الگوهای متناوب ارزش‌گذاری شده، شروع به یادگیری عملکرد سیستم می‌کند تا بتواند با استفاده از

شرایط حال حاضر سیستم، تخمین مناسبی از شرایطی که در آینده پیش‌خواهد آمد داشته باشد. در این رهیافت برای یادگیری ماشین از روش جنگل تصادفی^{۱۲}، استفاده شده است. در نهایت نیز برای بهبود بیشتر عملکرد، یک معیار اعتماد برای خروجی‌های تخمینی الگوریتم، تعریف می‌شود.

با توجه به این‌که این روش از یادگیری ماشین استفاده می‌کند و از داده‌های تاریخی استفاده می‌شود، این سیستم عملکرد مناسبی دارد. زیرا با تحلیل گذشته سیستم، به تخمین درستی از اتفاقات پیش‌رو می‌رسد و به همین خاطر می‌تواند بهترین تصمیم‌گیری‌ها در مورد نحوه تامین منابع را انجام دهد. مشکلی که در این روش وجود دارد، نیازمندی این روش به جمع‌آوری حجم زیادی از داده و صرف وقت نسبتاً زیاد برای یادگیری است. همچنین، این روش از جنگل تصادفی برای آموزش استفاده می‌کند که دارای پیچیدگی زیادی بوده و به قدرت محاسباتی بالایی نیاز دارد. با توجه به دو مورد گفته شده، در استفاده از این روش باید با بررسی سیستم به این نتیجه رسید که جمع‌آوری حجم زیاد داده و صرف وقت و قدرت محاسباتی بالا برای آموزش مدل، توجیه منطقی دارد.

در روش‌های مقیاس‌پذیری مبتنی بر یادگیری، الگوریتم مقیاس‌پذیری همواره در حال جمع‌آوری داده در مورد عوامل موثر در عملکرد و وضعیت سامانه از نظر مصرف منابع و خدمات ارائه شده به کاربر است. پس از جمع‌آوری حجم مناسبی از داده، الگوریتم مقیاس‌پذیری شروع به یادگیری شرایط سامانه می‌کند به این معنی که الگوریتم بتواند با داشتن متغیرهای سامانه مانند، بار ورودی و شرایط منابع، تخمینی از شرایطی که در پیش روی محیط رایانش ابری است، داشته باشد و با استفاده از این تخمین، قبل از رسیدن به نیاز به افزایش یا کاهش منابع، دستور انجام عملیات مناسب صادر شود. داده‌های جمع‌آوری شده بسته به شرایط هر مسئله متفاوت خواهد بود. برای پیاده‌سازی الگوریتم یادگیری نیز از روش‌های متفاوتی

10- TPM

11- Autocorrelation

12- Random Forest

جدول ۱: مقایسه کارهای انجام شده

روش ارائه شده	نویسندگان	نوع	سیاست	تکنیک
Auto-scaling model for cloud computing system[9]	Hung, Che-Lun and Hu, Yu-Chen and Li, Kuan-Ching	واکنشی	عمودی	مبتنی بر آستانه
Auto-scaling Strategy for Amazon Web Services in Cloud Computing[14]	Liao, Wen-Hwa and Kuai, Ssu-Chi and Leau, Yu-Ren	واکنشی	عمودی	مبتنی بر آستانه
A Rapid Auto-Scaling Mechanism in Cloud Computing Environment[11]	Tseng, Chia-Wei and Tsai, Ming-Shiun and Yang, Yao-Tsung and Chou, Li-Der	واکنشی	عمودی	مبتنی بر استنباط فازی
An auto-scaling mechanism for cloud-based multimedia storage systems: A fuzzy-based elastic controller[8]	Mostafa Ghobaei-Arani, Maryam Rezaei, Alireza Souri	واکنشی	عمودی	مبتنی بر استنباط فازی
AME-WPC: Advanced model for efficient workload prediction in the cloud[10]	Katja Cetinski and Matjaz B. Juric	فعال	عمودی	مبتنی بر تطبیق الگو
A reinforcement learning based auto-scaling approach for SaaS providers in dynamic cloud environment[12]	Wei, Yi and Kudenko, Daniel and Liu, Shijun and Pan, Li and Wu, Lei and Meng, Xiangxu	فعال	عمودی	مبتنی بر تطبیق الگو
A Comparison of Reinforcement Learning Techniques for Fuzzy Cloud Auto-Scaling[13]	Arabnejad, Hamid and Pahl, Claus and Jamshidi, Pooyan and Estrada, Giovanni	ترکیبی	عمودی	مبتنی بر استنباط فازی / یادگیری

به دنبال پیدا کردن حالت بهینه اختصاص منابع بوده و برای رسیدن به این هدف از استنباط فازی استفاده می‌کند. همچنین، برای مدیریت بهتر از مدل حلقه OODA در این روش استفاده شده است. رهیافت پیشنهادی از عواملی مانند میزان مصرف پردازنده و مدت زمان پاسخگویی محیط رایانش ابری برای تصمیم‌گیری در مورد نحوه تامین منابع استفاده می‌کند. تلاش رهیافت پیشنهادی این است که بتواند با تعیین مقدار مناسب منابع مورد نیاز، همزمان با انجام دادن درخواست‌های کاربران، از صرف هزینه‌های اضافه خودداری کند. روش پیشنهاد شده در این کار تحقیقاتی یک روش مقیاس‌پذیری پویا خودکار مبتنی بر منطق و استنباط فازی است که تلاش می‌کند متناسب با شرایط محیط رایانش ابری، تصمیمات و عملیات مقیاس‌پذیری را انجام دهد. محیط رایانش ابری در رهیافت پیشنهاد شده شامل قسمت‌های متفاوتی مانند، مرکز داده^{۱۴}، ماشین‌های فیزیکی^{۱۵}، ماشین‌های مجازی، واسط^{۱۶} و کارهای درخواستی از طرف کاربران^{۱۷} است. معماری کلی محیط رایانش ابری در شکل ۲ نشان داده شده است.

مانند Q-Learning، یادگیری تقویتی^{۱۳} [۱۲]، یادگیری ماشین و شبکه‌های عصبی [۶] استفاده می‌شود. روش‌های مقیاس‌پذیری ترکیبی از ترکیب چندین روش با یکدیگر تولید می‌شوند. به‌طور مثال با ترکیب روش Q-Learning و منطق فازی می‌توان یک روش جدید برای مقیاس‌پذیری خودکار ایجاد کرد. نقطه مثبت روش‌های ترکیبی این است که در صورتی که روش‌های مناسب با یکدیگر ترکیب شوند، می‌توانند مکمل یکدیگر باشند و هر روش ایرادات روش دیگر را جبران کند. به‌طور مثال در کار تحقیقاتی انجام شده توسط خانم عرب‌نژاد و همکاران ایشان [۱۳]، یک روش fuzzy Q-Learning معرفی شده است. در این روش سعی شده است تا با ترکیب منطق فازی با یادگیری تقویتی، از طرفی مشکل ساده‌بودن و عدم استفاده از داده‌های محیط در روش فازی و از طرف دیگر، مشکل نیاز به ذخیره‌سازی حالات و اقدامات در روش‌های یادگیری تقویتی، مرتفع شود.

۴- رهیافت پیشنهادی

در رهیافت ارائه شده، الگوریتم مقیاس‌پذیری خودکار

13- Reinforcement Learning

14- Data Center
15- Host
16- broker
17- Cloudlet

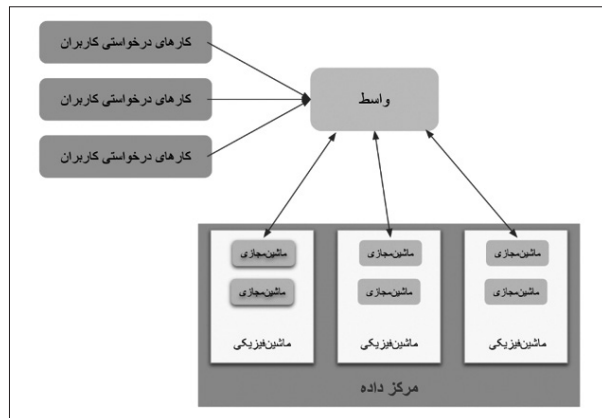
میانگین پردازنده مصرفی به عنوان یکی از مؤلفه‌های موثر در مقیاس‌پذیری خودکار در نظر گرفته شده است. نحوه محاسبه این مؤلفه در فرمول ۱ قابل مشاهده است که در آن MCU میانگین پردازنده مصرف شده توسط تمامی ماشین‌های مجازی، CPU_{VM_i} مقدار پردازنده مصرف شده توسط ماشین مجازی i ام و $|VM|$ تعداد ماشین‌های مجازی موجود در مرکز داده است.

$$MeanCpuUsagt = MCU = \frac{\sum CPU_{VM_i}}{|VM|} \times 100 \quad (1)$$

۴-۱-۲- تاخیر پاسخگویی به درخواست کاربران

منابع محاسباتی موجود در مرکز داده برای انجام درخواست‌های کاربران نیاز به مدت زمانی برای اجرای کامل درخواست‌ها دارند. مشخص است که هر چه مدت زمان اجرای درخواست‌ها کمتر باشد، عملکرد سیستم بهتر خواهد بود. نکته‌ای که قابل توجه است این است که با استفاده از برخی روش‌ها می‌توان زمان اجرای درخواست کاربران را کاهش داد اما با توجه به قدرت محاسباتی و حجم کار ورودی، مدت زمان پاسخگویی به درخواست کاربر، دارای یک کمینه بوده که مدت زمان پاسخگویی از آن کمینه، کمتر نمی‌شود. در مسئله مورد بررسی، کارهای درخواستی کاربران، دارای مهلت اجرای سخت^{۱۹} نبوده که در صورتی که سیستم نتواند آن را رعایت کند دچار مشکل شود. اما با توجه به حجم درخواست کاربران و قدرت محاسباتی منبع، کمترین زمانی که یک منبع محاسباتی ایده‌آل می‌تواند برای اجرای درخواست نیاز داشته باشد، به عنوان مهلت اجرای نرم^{۲۰}، مشخص می‌شود. هرچه مدت زمان اجرای درخواست کاربر از این حد کمینه بیشتر شود، نشان دهنده پایین آمدن عملکرد سیستم و در نتیجه احتمال وجود کمبود منابع می‌باشد. همچنین در شرایطی که هر درخواست در زمان کمینه خود اجرا شود، می‌توان این موضوع را برداشت کرد که سیستم در یک حالت ایده‌آل برای اجرای درخواست‌ها قرار دارد. برای محاسبه این مؤلفه مطابق با فرمول ۲ عمل می‌شود. در این فرمول

19- Hard deadline
20- Soft deadline



شکل ۲: معماری کلی محیط رایانش ابری

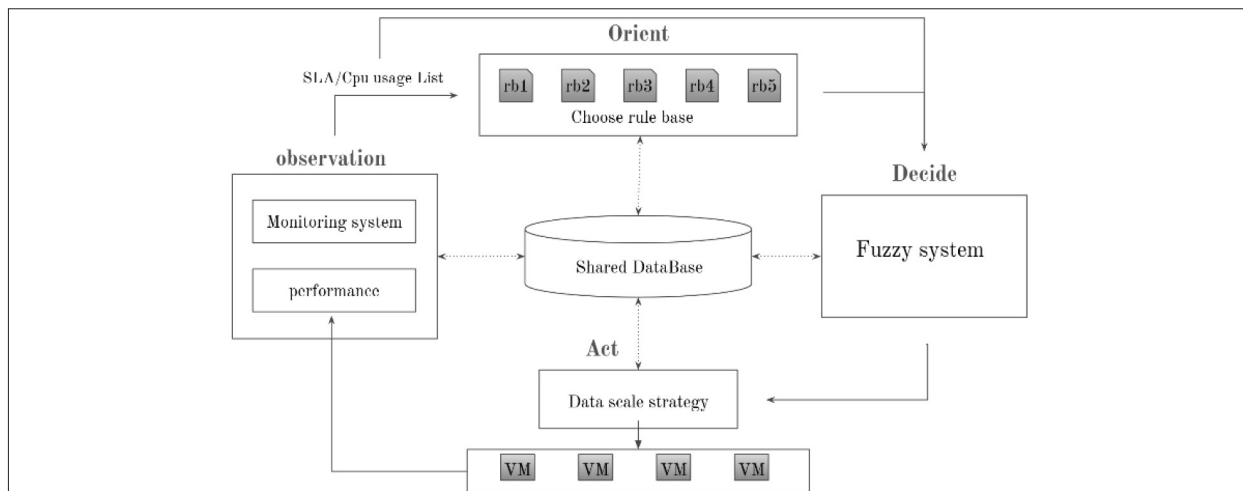
۴-۱- مؤلفه‌های مورد بررسی در رهیافت پیشنهادی

همان‌طور که گفته شد، متغیرها و مؤلفه‌های متفاوتی در تصمیم‌گیری مقیاس‌پذیری درباره تامین منابع موثر هستند و متناسب با هر مسئله و شرایط خاص آن مسئله متفاوت خواهند بود. در این مسئله، دو مؤلفه میانگین مقدار پردازنده مصرف شده و تاخیر پاسخگویی به درخواست‌های کاربران، به عنوان دو مؤلفه مورد بررسی و موثر در تصمیم‌گیری انتخاب شده‌اند که در ادامه به توضیح هر کدام می‌پردازیم.

۴-۱-۱- میانگین مقدار پردازنده مصرف شده^{۱۸}

هر ماشین فیزیکی و ماشین مجازی در مرکز داده، دارای تعدادی هسته پردازشی و در نتیجه قدرت محاسباتی خاص خود هستند. هنگامی که کارهای درخواستی کاربران برای انجام به ماشین‌های مجازی سپرده می‌شود، هر ماشین مجازی با استفاده از منابع محاسباتی خود مانند پردازنده، تلاش می‌کند که در کمترین زمان ممکن وظیفه خود را انجام دهد. هرچه بارکاری ورودی به مرکز داده و در نتیجه ماشین‌های فیزیکی و مجازی بیشتر باشد، نیاز به قدرت محاسباتی بیشتر بوده و در صورتی که تعداد منابع محاسباتی یا ماشین‌های مجازی ثابت باشند، هر ماشین مجازی با مصرف قدرت محاسباتی و مصرف پردازنده بیشتر، کارهای ورودی را کنترل می‌کند. با توجه به موضوع گفته شده، در رهیافت پیشنهادی مقدار

18- CPU



شکل ۳: معماری مدل حلقه OODA در رهیافت پیشنهادی

رهیافت پیشنهادی در آن پیاده‌سازی می‌شود.

الگوریتم ۱: شبه کد الگوریتم مرحله Observation

شبه کد الگوریتم مرحله Observation

Algorithm: Observation Phase
 Input: List CPU, List ResponseTime, List DeadLine
 Output: SLA, MCU
 MCU = 0
 SLA = 0
 | For i = 1 to |VM
 (MCU += CPU.get(i)
 (If ResponseTime.get(i) > DeadLine.get(i)
 ++ SLA
 End if
 End for
 | MCU = (MCU * 100) / |VM
 | SLA = (SLA * 100) / |VM
 Return MCU, SLA

وظیفه اصلی این مرحله، بررسی شرایط محیط با استفاده از مؤلفه‌های ضبط شده در پایگاه داده در مرحله قبل و سپس تصمیم‌گیری در مورد نحوه انتخاب پایگاه قاعده برای استفاده در استتباط فازی در مراحل بعدی، است. ابتدا، الگوریتم حالت سامانه را از میان پنج حالت از پیش تعریف شده انتخاب می‌کند. این پنج حالت شامل تامین اضافی منابع در حالت شدید، تامین اضافی منابع، حالت عادی، کمبود در منابع و کمبود بیش از اندازه در منابع، تعریف شده‌اند. پس از تشخیص حالت سامانه، پایگاه قاعده متناسب با شرایط حال حاضر انتخاب می‌شود. انتخاب پایگاه قاعده با استفاده از دو مؤلفه مورد بررسی، MCU و SLA، به صورت زیر انجام می‌شود. حالت‌های ExOver

SLAT، Responsetime، DeadLine و |UserRequests| به ترتیب نشان‌دهنده نقض شدن یا نشدن توافقنامه سطح خدمات برای درخواست آام، مدت زمان اجرای درخواست آام، مهلت تعیین شده اجرای درخواست آام و تعداد درخواست‌های ورودی به مرکز داده هستند.

$$SLA_T = \begin{cases} 1 & \text{Responsetime} > \text{DeadLine} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

$$SLA = \frac{\sum SLA_T}{|UserRequests|} \times 100$$

۴-۲ الگوریتم روش پیشنهادی

همان‌طور که گفته شد، در رهیافت پیشنهادی از مدل حلقه OODA که شامل چهار مرحله، مشاهده، جهت‌گیری، تصمیم‌گیری و عملیات است استفاده می‌شود. معماری کلی رهیافت پیشنهادی با توجه به حلقه OODA در شکل ۳ قابل مشاهده است.

۴-۲-۱ مرحله Observation

این مرحله، اولین مرحله اجرای الگوریتم مقیاس‌پذیری خودکار می‌باشد. در این مرحله اطلاعات اولیه درباره شرایط محیط رایانش ابری، برای استفاده در مراحل بعدی، جمع‌آوری می‌شود و در یک پایگاه داده مشترک ذخیره می‌شود.

۴-۲-۲ مرحله Orientation

این مرحله یکی از قسمت‌هایی است که پویایی مضاعف

متناسب با پایگاه قاعده و شرایط محیط، متفاوت است. پس از انجام عملیات فازی سازی متغیرها، از پایگاه قاعده برای ترکیب قوانین موجود در پایگاه قاعده با متغیرهای ورودی فازی سازی شده، استفاده می‌شود. برای ترکیب ورودی با قوانین باید قدرت و مقدار درستی هر قانون با توجه به مقادیر ورودی به دست آید که در این مسئله، با توجه به این که از منطق استنباط فازی به روش ممدانی استفاده شده است، از روش کمینه برای تعیین درجه درستی قوانین استفاده شده است. پس از محاسبه مقادیر خروجی برای هر قانون باید این مقادیر با یکدیگر ترکیب و جمع شوند و برای این کار از بیشینه استفاده می‌شود. پس از ترکیب خروجی‌ها، نیاز به خنثی‌سازی متغیر خروجی که همان مؤلفه مقیاس‌گذار می‌باشد، وجود دارد. با توجه به استفاده از روش ممدانی در این مرحله، خروجی، مجموعه‌ای از سطوح زیر نمودارهای فازی برای متغیر مقیاس‌پذیری در قوانین می‌باشد، در نتیجه استفاده از روش مرکز ثقل، منطقی‌ترین روش برای خنثی‌سازی متغیر خروجی می‌باشد. در مرحله آخر برای تصمیم‌گیری نهایی در رابطه با نحوه مقیاس‌پذیری، با توجه به مقدار عددی متغیر مقیاس‌پذیری (scale) و فرمول ۴، دستور مقیاس‌پذیری در پایگاه داده مشترک ذخیره می‌شود. متغیر Increase به معنای افزایش در منابع، Decrease به معنای کاهش در منابع و NoChange به معنای عدم تغییر در منابع است.

$$scale = \begin{cases} Increase & scale \geq 60 \\ NoChange & 50 < scale < 60 \\ Decrease & scale \leq 50 \end{cases} \quad (4)$$

الگوریتم ۳: شبه کد الگوریتم مرحله Decision

شبه کد الگوریتم مرحله Decision
Algorithm: Decision Phase
Input: RuleBaseIndex, SLA, MCU
Output: scale
CurrentRuleBase = loadFromRuleBaseCollection(RuleBaseIndex)
(FuzzyInference f = new FuzzyInference(CurrentRuleBase))
(f.setInputVariable(SLA, MCU))
(scale = f.inference)
Return scale

ExUnder, Under, over, و Normal به ترتیب متناظر با حالت‌های تامین اضافی منابع در حالت شدید، تامین اضافی منابع، کمبود در منابع، کمبود بیش از اندازه در منابع و حالت عادی هستند.

$$RuleBase = \begin{cases} ExOver & SLA = 0 \text{ and } MCU \leq 15 \\ Over & SLA = 0 \text{ and } 15 < MCU \leq 30 \\ Under & 50 \leq SLA \leq 75 \text{ and } 30 < MCU \leq 60 \\ ExUnder & 75 < SLA \text{ and } 70 \leq MCU \\ Normal & else \end{cases} \quad (3)$$

در ادامه الگوریتم این مرحله از روش پیشنهادی

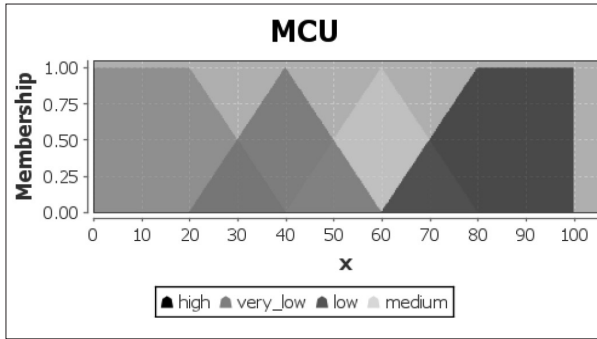
مشاهده می‌شود
الگوریتم ۲: شبه کد الگوریتم مرحله Orientation

شبه کد الگوریتم مرحله Orientation

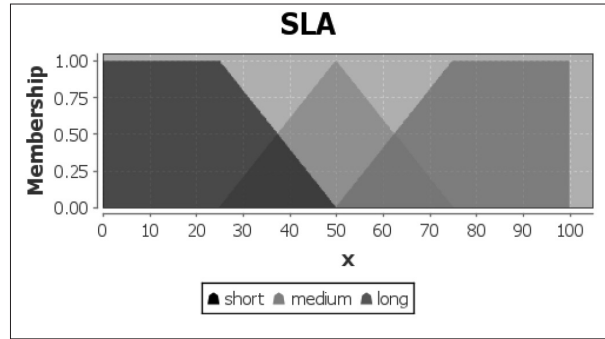
```
Algorithm: Orientation Phase
Input: SLA, MCU
Output: RuleBaseIndex
RuleBaseIndex = 3
If SLA = 0
If MCU <= 15
RuleBaseIndex = 1
End if
Else
RuleBaseIndex = 2
End else
End if
Else
If 30 <= MCU <= 60 or 50 <= SLA <= 75
RuleBaseIndex = 4
End if
Else if 70 <= MCU or 75 < SLA
RuleBaseIndex = 5
End else if
End else
Return RuleBaseIndex
```

۴-۲-۳- مرحله Decision

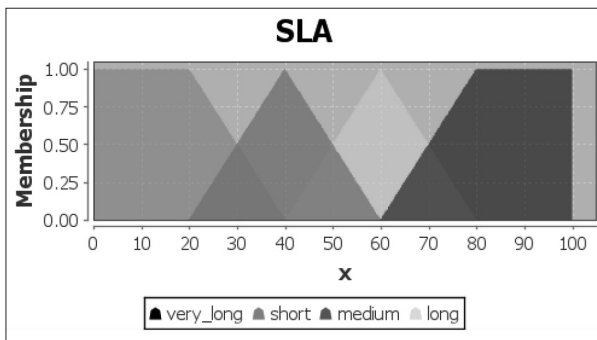
پس از انتخاب پایگاه قاعده متناسب با شرایط محیط رایانش ابری در مرحله قبل و جمع‌آوری مؤلفه‌های موثر در سیستم، در مرحله مشاهده، الگوریتم پیشنهادی در این مرحله به استنباط فازی برای به دست آوردن بهترین عملیات مقیاس‌پذیری می‌پردازد. پس از دریافت مقدار مؤلفه‌ها، اولین مرحله، فازی‌سازی دو مؤلفه SLA و MCU است. برای فازی‌سازی مؤلفه‌ها نیاز است تا اصطلاحات فازی هر کدام و تابع فازی‌سازی آن‌ها تعیین شود. نحوه فازی‌سازی هر دو متغیر مطابق با شکل‌های ۴ تا ۹ است. همان‌طور که مشخص است، برای فازی‌سازی مؤلفه‌ها از تابع فازی‌سازی مثلثی استفاده شده است و فازی‌سازی



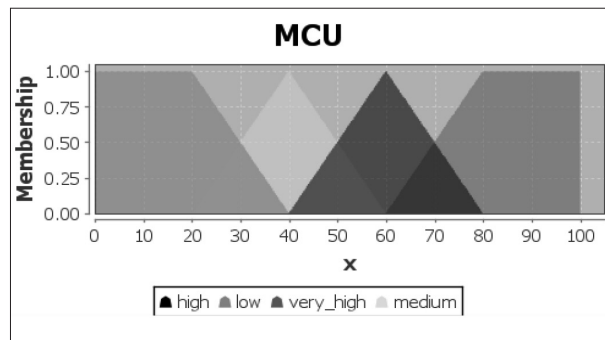
شکل ۵: فازی سازی مؤلفه MCU



شکل ۴: فازی سازی مؤلفه SLA



شکل ۷: فازی سازی مؤلفه SLA در حالت کمبود منابع شدید



شکل ۶: فازی سازی مؤلفه MCU در حالت کمبود منابع شدید

تأمین اضافی منابع بیش از حد باشد، دو ماشین مجازی حذف می‌شوند و در حالت‌های دیگر یک ماشین مجازی حذف می‌شود.

الگوریتم ۴: شبه کد الگوریتم مرحله Action

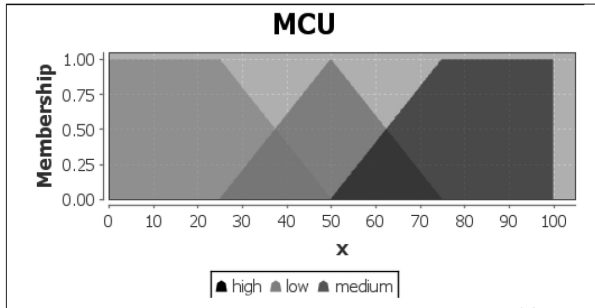
شبه کد الگوریتم مرحله Action
Algorithm: Action Phase
Input: scale, RuleBaseIndex, List ExecutingVms
Output: scaling action
step = 1
If RuleBaseIndex = 5 or RuleBaseIndex = 1
step = 2
End if
If 60 <= scale
For i = 1 to step
(AddNewVM
End for
End if
Else if scale <= 50
;(idleVm = ExecutingVms.get(0
For i = 1 to ExecutingVms
If ExecutingVms.get(i).IsIdle
(idleVm = ExecutingVms.get(i
End if
End for
(RemoveVM(idleVm
End Else id

۴-۲-۴- مرحله Action

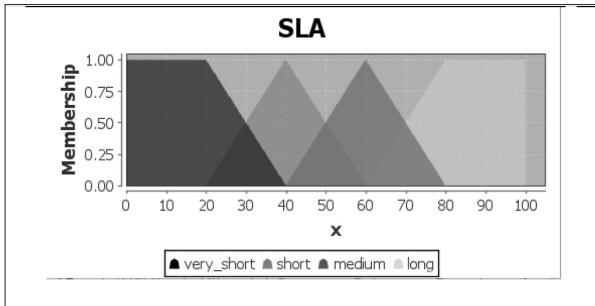
قدم آخر در الگوریتم پیشنهادی، پیاده‌سازی دستور صادر شده در مرحله قبل توسط مؤلفه مقیاس‌پذیری خودکار است. این مرحله نیز یکی دیگر از مراحل است که وظیفه پیاده‌سازی قسمتی از پویایی اضافه شده به روش‌های فازی را بر عهده دارد. این مرحله در صورت صدور دستور افزایش در منابع موجود در مرکز داده، درخواست ایجاد ماشین مجازی را به واسط ارسال می‌کند. همان‌طور که گفته شد، تعداد ماشین‌های مجازی اضافه شده وابسته به شرایط محیط و پایگاه قاعده انتخاب شده، هستند. اگر پایگاه قاعده انتخاب شده متناسب با شرایط کمبود شدید منابع باشد، الگوریتم در این مرحله درخواست اضافه شدن دو ماشین مجازی به مرکز داده موجود را به واسط ارسال می‌کند و در حالت‌های دیگر درخواست اضافه شدن تنها یک ماشین مجازی را ارسال می‌کند. در دستور کاهش نیز پویایی در قدم مقیاس‌پذیری وجود دارد به این صورت که اگر پایگاه قاعده متناسب با شرایط

۵- ارزیابی رهیافت پیشنهادی

باتوجه به مشکلات پیاده‌سازی الگوریتم در محیط‌های واقعی رایانش ابری، در این کار پژوهشی از شبیه‌سازی برای ارزیابی روش پیشنهادی استفاده شده است. ارزیابی الگوریتم در این بخش در سه مرحله انجام می‌شود. مرحله اول تحلیل حساسیت رهیافت پیشنهادی است، به این صورت که با تغییر مقادیر مختلف در الگوریتم، نتیجه و عملکرد رهیافت پیشنهادی چگونه تغییر می‌کند. مرحله دوم، بررسی تاثیر مؤلفه‌های انتخابی ارزیابی است و مرحله سوم مقایسه رهیافت پیشنهادی با روش‌های پیشنهادی دیگر خواهد بود. در ادامه، ابتدا به معرفی ابزار شبیه‌ساز مورد استفاده و محیط آزمایش می‌پردازیم. سپس به ارزیابی رهیافت پیشنهادی مقیاس‌پذیری خودکار در سه مرحله گفته شده پرداخته می‌شود.



شکل ۸: فازی سازی مؤلفه MCU در حالت تامین اضافی منابع شدید



شکل ۹: فازی سازی مؤلفه SLA در حالت تامین اضافی منابع شدید

مرکز داده و تعدادی میزبان یا ماشین‌فیزیکی است. هر ماشین‌فیزیکی دارای منابع محاسباتی بوده که در اختیار ماشین‌های مجازی برای اجرای درخواست کاربران قرار می‌گیرد. همچنین برای پیاده‌سازی منطق فازی در روش پیشنهادی، از کتابخانه jFuzzyLogic استفاده شده است. محیط انجام آزمایش شامل یک مرکز داده و تعدادی میزبان یا ماشین‌فیزیکی است. هر ماشین‌فیزیکی دارای منابع محاسباتی است که در اختیار ماشین‌های مجازی برای اجرای درخواست کاربران قرار می‌گیرد. در مسئله‌های مختلف، بار ورودی و درخواست‌های کاربران می‌تواند از الگوی خاصی پیروی کند و یا این‌که به صورت تصادفی باشد. در ارزیابی انجام شده، بار ورودی توسط کاربران دارای هیچ الگوی مشخصی نبوده و به صورت تصادفی، درخواست‌های کاربران به واسطه ارسال می‌شود. هر درخواست دارای اندازه مشخصی است که نشان دهنده تعداد دستورات مورد نیاز برای انجام این درخواست است. در ارزیابی انجام شده، طول و اندازه درخواست‌های کاربران به صورت تصادفی انتخاب می‌شوند. نحوه استفاده از منابع محاسباتی ماشین‌های مجازی اختصاص داده شده به درخواست برای اجرا، که در ارزیابی انجام شده یک درخواست به طور کامل از

الگوریتم ۵: شبه کد الگوریتم مقیاس‌پذیری خودکار پیشنهادی

شبه کد الگوریتم مقیاس‌پذیری خودکار پیشنهادی

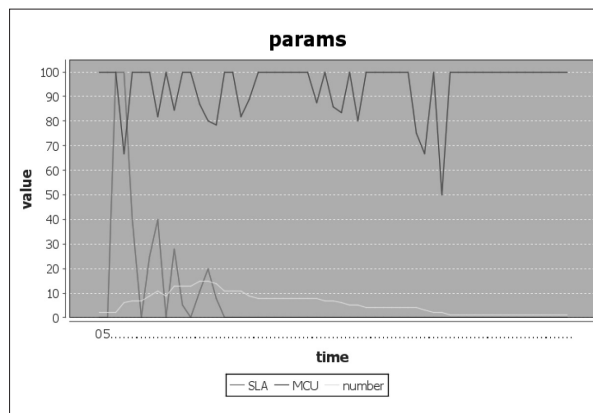
```
Algorithm: ProposedAutoscalingApproach
Input: MCU, SLA
Output: scale
While (TRUE) do
Begin While
|For i = 1 to |requests
(MainQueue.Add(requesti
End for
(Sorted_Req = Priority(MainQueue,DeadLine
(Observation.ADD(DB.Request,sorted_Req
(Observation.ADD(DB.VM,VMs
(requesti,VMs)=Orientation.Read(DB.Request,DB.VM)
(Rulebase = choose_rule_base(requests,VMs
(Orientation.ADD(DB.Rulebase,Rulebase
(requests,VMs,rulebase) = Decision.Read(DB.Requests,DB.)
(VMs,DB.Rulebase
Strategy = Scale_Strategy_Fuzzy_
(logic(Requests,VMs,rulebase
(Decision.ADD(DB.Strategy,Strategy
(Req,command) = Action.Read(DB.Request,DB.Strategy)
(Action(Req,command
End While
```

طراحی و پیاده‌سازی رهیافت پیشنهادی مبتنی بر استنباط فازی با استفاده از زبان جاوا و در محیط شبیه‌سازی Cloud-Sim Plus انجام شده است. محیط انجام آزمایش شامل یک

در شکل ۱۰ قابل مشاهده است.

الگوریتم در ابتدای اجرا از نظر مؤلفه تأخیر در پاسخگویی عملکرد مناسبی ندارد و مقدار این مؤلفه تقریباً صد درصد می‌باشد اما با گذشت زمان و با اضافه کردن ماشین‌های مجازی بیشتر به مرکز داده، رفته رفته مقدار این تأخیر کاهش پیدا می‌کند. همچنین از نظر مؤلفه مقدار CPU مصرف شده، الگوریتم تلاش می‌کند تا مقدار CPU را در بیشینه خود، یعنی صد درصد، قرار ندهد و این موضوع نیز در نمودار قابل مشاهده است. از نظر تعداد ماشین‌های مجازی نیز الگوریتم پس از این‌که در مؤلفه تأخیر پاسخگویی عملکرد نامناسبی پیدا می‌کند، تصمیم به افزایش تعداد ماشین‌های مجازی می‌گیرد و پس از کاهش این تأخیر، تعداد ماشین‌های مجازی نیز کاهش پیدا می‌کند. نکته قابل توجه در خروجی رهیافت این است که از آنجایی که روش پیشنهادی یک روش واکنشی می‌باشد، افزایش تعداد ماشین‌های مجازی با یک تأخیر انجام می‌شود به این معنی که هنگامی که الگوریتم با توجه به متغیر SLA متوجه تأخیر زیاد در پاسخگویی به درخواست‌های کاربران می‌شود، دستور به افزایش در تعداد ماشین‌های مجازی می‌دهد اما همان‌طور که در نمودار قابل توجه است این افزایش با مقداری تأخیر نسبت به نقطه اوج متغیر SLA انجام می‌شود.

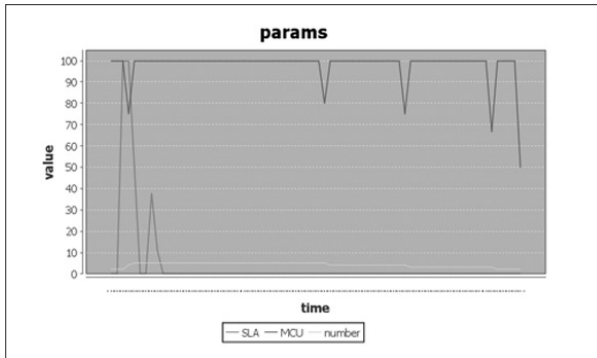
برای بررسی تأثیر پویایی‌های اضافه شده به الگوریتم، در دو مرحله پویایی اضافه شده به الگوریتم در انتخاب پایگاه قاعده و در انتخاب قدم مقیاس‌پذیری را حذف می‌کنیم تا نتایج آن را بررسی کنیم. برای ایجاد ایستایی در پایگاه قاعده، به جای انتخاب یک پایگاه قاعده از میان پنج پایگاه قاعده از پیش تعریف شده، پایگاه قاعده متناسب با شرایط عادی، که در تصمیم‌گیری هیچ سوگیری در افزایش یا کاهش منابع ندارد، به عنوان پایگاه قاعده استفاده می‌شود. نتایج این روش در شکل ۱۲ قابل مشاهده است. همان‌طور که مشخص است، در روش ایستا، مؤلفه SLA نسبت به حالت پویا، در بازه‌های زمانی یکسان، دارای مقادیر بیشتری است. همچنین مقدار CPU مصرفی در مرکز داده نسبت به حالت پویا، مدت



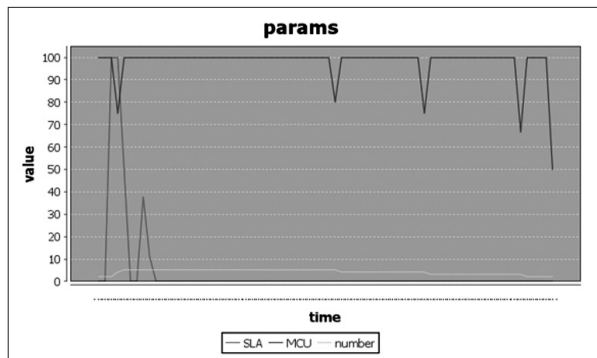
شکل ۱۰: خروجی رهیافت پیشنهادی

منابع محاسباتی ماشین‌های مجازی استفاده می‌کند. در ارزیابی انجام شده، محیط رایانش ابری شامل تنها یک مرکز داده با تعدادی ماشین‌فیزیکی یا میزبان است. همچنین هر میزبان دارای تعدادی هسته محاسباتی بوده که توانایی هر هسته اجرای ۱۰۰۰ میلیون دستور در ثانیه است. هر ماشین‌فیزیکی تعدادی ماشین‌های مجازی را میزبانی می‌کند. وظیفه اجرای درخواست‌های کاربران بر عهده ماشین‌های مجازی است. هر ماشین‌های مجازی منابع خود را از ماشین‌فیزیکی خود تأمین می‌کند و در نتیجه همانند ماشین‌های فیزیکی، یک ماشین‌های مجازی دارای ویژگی‌های محاسباتی خواهد بود. هر ماشین‌های مجازی دارای یک ظرفیت برای اجرای درخواست‌های کاربران بوده که ظرفیت تمامی ماشین‌های مجازی موجود در مرکز داده ۱۰۰۰ میلیون دستور در ثانیه است. همچنین هر ماشین‌های مجازی دارای دو هسته محاسباتی است.

در روش پیشنهاد شده، الگوریتم سعی دارد با کم کردن مقدار پردازنده مصرف شده توسط منابع محاسباتی و همچنین کمینه کردن مدت پاسخگویی به درخواست کاربران، به بهترین تصمیم در مسئله مقیاس‌پذیری در محیط رایانش ابری دست پیدا کند. همچنین، تعداد ماشین‌های مجازی استفاده شده برای انجام درخواست‌های کاربران نیز عامل موثری در ارزیابی کارایی روش پیشنهادی است. در نتیجه، سه مؤلفه SLA، MCU و تعداد ماشین‌های مجازی موجود در مرکز داده به عنوان شاخص‌های ارزیابی روش پیشنهادی در نظر گرفته شده‌اند. خروجی الگوریتم بر اساس این سه مؤلفه

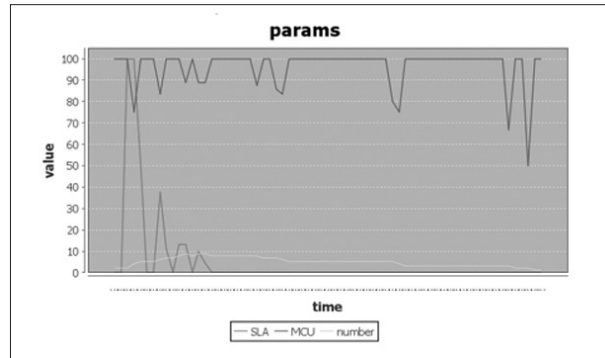


شکل ۱۲: خروجی رهیافت پیشنهادی بدون پایگاه قاعده پویا



شکل ۱۳: خروجی روش مبتنی بر استنباط فازی

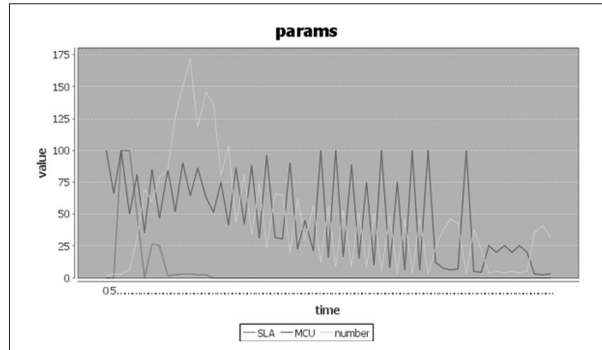
الگوریتم در قدم‌های مقیاس‌پذیری واضح است. هر چند که این تاثیر به اندازه تاثیر مثبت پایگاه قاعده پویا نیست. برای ارزیابی بهتر روش ارائه شده، رهیافت پیشنهادی با دو روش دیگر نیز مقایسه شده است. در قسمت اول به مقایسه این روش با یک روش مقیاس‌پذیری استنباط فازی می‌پردازیم. روش مقیاس‌پذیری خودکار مبتنی بر استنباط فازی، یک روش مشابه با رهیافت پیشنهادی است با این تفاوت که در این حالت الگوریتم تماماً ایستا بوده و تنها از یک پایگاه قاعده و یک قدم مقیاس‌پذیری در آن استفاده می‌شود. برای ارزیابی روش مورد مقایسه، پایگاه قاعده همان پایگاه قاعده حالت عادی مورد استفاده در رهیافت پیشنهادی است و قدم مقیاس‌پذیری، یک ماشین مجازی است. خروجی روش مورد مقایسه و رهیافت پیشنهادی در شکل ۱۳ قابل مشاهده است. با مقایسه خروجی روش مورد مقایسه با روش پیشنهادی، مشخص است که مؤلفه MCU در روش استنباط فازی نسبت به رهیافت پیشنهادی، دارای مقادیر بیشتری بوده و در زمان‌های بیشتری در حالت بیشینه است.



شکل ۱۴: خروجی رهیافت پیشنهادی بدون قدم مقیاس‌پذیری پویا

زمان بیشتری در حالت بیشینه خود، یعنی صد درصد خواهد بود. از نظر تعداد ماشین‌های مجازی نیز الگوریتم در حالت ایستا، تغییرات کمتری نسبت به حالت پویا دارد و به نظر می‌آید متناسب با شرایط تغییر نمی‌کند. با توجه به موارد گفته شده می‌توان نتیجه گرفت پویایی اضافه شده در مرحله انتخاب پایگاه قاعده، مؤلفه موثری در بهبود عملکرد رهیافت پیشنهادی است.

برای بررسی تاثیر قدم مقیاس‌پذیری پویا و ایجاد ایستایی در آن، به جای استفاده از تعداد قدم متغیر، در افزایش و یا کاهش منابع، تنها یک ماشین مجازی اضافه و یا کم می‌شود. نتایج خروجی این روش در شکل قابل مشاهده است. با توجه به خروجی الگوریتم در حالت ایستای قدم‌های مقیاس‌پذیری، از نظر مقدار CPU مصرفی این حالت نیز مانند حالت قبل، مدت زمان بیشتری در حالت بیشینه خود است اما نسبت به حالت قبلی، این مدت زمان کمتر خواهد بود. این شرایط برای مؤلفه SLA نیز برقرار بوده به این معنی که حالت ایستای قدم‌های مقیاس‌پذیری نسبت به حالت پویا عملکرد نامناسب تری دارد. اما نسبت به حالت ایستای پایگاه قاعده، عملکرد آن بهتر است. علت عملکرد بهتر این روش نسبت به حالت ایستای پایگاه قاعده این است که تغییر ایجاد شده در این حالت بسیار جزئی بوده و تنها در سرعت عمل الگوریتم موثر است. از طرفی علت عملکرد بهتر روش پویا نسبت به این روش، به دلیل امکان ایجاد تغییرات گسترده‌تر در حالت‌هایی که وضعیت سیستم اضطرابی است خواهد بود. با توجه به موضوعات بررسی شده، تاثیر پویایی اضافه شده به



شکل ۱۴: خروجی روش مبتنی بر آستانه

در شکل ۱۴ نمایش داده شده است. همان طور که در شکل قابل مشاهده است، با توجه به این که این روش یک آستانه بالایی برای مقدار میانگین CPU مصرفی در نظر گرفته است، مقدار میانگین CPU در این روش نسبت به روش پیشنهادی کم تر بوده اما این روش در کاهش ماشین های مجازی بسیار محتاط عمل می کند. به همین دلیل تعداد ماشین های مجازی مورد استفاده در این روش نسبت به روش پیشنهادی بسیار بیشتر است و می توان گفت که در این روش، محیط رایانش دچار تامین بیش از اندازه منابع محاسباتی شده است. از نظر تاخیر در پاسخگویی به درخواست کاربران نیز با این که این روش نسبت به رهیافت پیشنهادی از تعداد بیشتری ماشین مجازی بهره می برد، اما تاخیر پاسخگویی در این روش تفاوت چندانی با رهیافت پیشنهادی ندارد.

با توجه به موارد گفته شده می توان نتیجه گرفت این روش از نظر مقدار CPU مورد استفاده در محیط رایانش ابری، نسبت به رهیافت پیشنهادی عملکرد بهتری دارد اما باید توجه داشت این موضوع با افزایش تعداد ماشین های مجازی و در نتیجه افزایش هزینه های تامین این ماشین های مجازی همراه است و از طرفی، کم شدن مقدار مصرف CPU تاثیری چندانی در کاهش تاخیر پاسخگویی به درخواست کاربران نداشته است. به بیانی دیگر می توان گفت در این روش محیط رایانش ابری دچار تامین بیش از اندازه منابع محاسباتی شده است.

۶- نتیجه گیری و کارهای آینده

در رهیافت پیشنهاد شده تلاش بر این بوده است که با اضافه کردن پویایی به روش های مقیاس پذیری خودکار مبتنی بر استنباط فازی، عملکرد بهتری در زمینه مقیاس پذیری حاصل شود. برای ایجاد پویایی در الگوریتم تلاش شده است که ابتدا شرایط محیط تشخیص داده شود به این معنی که محیط رایانش ابری در حال حاضر در حالت کمبود در تامین منابع یا در حالت تامین بیش از حد منابع است. همچنین اگر به طور مثال محیط در حالت کمبود در تامین منابع باشد،

مؤلفه تاخیر نیز در این روش نسبت به روش پیشنهادی، دارای مقادیر بیشتری است. از نظر تعداد ماشین های مجازی مورد استفاده نیز با توجه به تغییرات بسیار کم در تعداد آن ها و با توجه به بالا بودن مقدار پردازنده مصرفی به نظر می آید که محیط رایانش ابری دچار کمبود در تامین منابع است. به طور کلی و با توجه به سه مؤلفه مورد بررسی برای ارزیابی الگوریتم، به نظر می آید که روش استنباط فازی نسبت به روش پیشنهاد شده دارای عملکرد ضعیف تری بوده و دچار کمبود در تامین منابع است. علت این موضوع این است که روش مورد مقایسه، هیچکدام از دو مرحله پویایی پیشنهاد شده، که هر دو باعث بهبود عملکرد الگوریتم می شوند، را ندارد و در نتیجه عملکرد ضعیف تری نسبت به رهیافت پیشنهادی از خود نشان می دهد.

روش دیگری که در ارزیابی مورد مقایسه قرار می گیرد یک روش مقیاس پذیری خودکار مبتنی بر آستانه و بر اساس روش ارائه شده توسط هانگ و همکارانشان [۹] است. این روش نیز در محیط شبیه سازی کاملاً مشابه با محیط ارزیابی رهیافت پیشنهادی، اجرا و مقایسه شده است. در این روش، در حالتی که مقدار میانگین CPU مصرفی از یک آستانه مشخص، که در اینجا ۷۰ درصد در نظر گرفته شده است، بیشتر باشد، درخواست اضافه شدن یک ماشین مجازی به واسط ارسال می شود. برای کاهش منابع نیز در صورتی که میانگین CPU مصرفی از ۳۰ درصد کم تر باشد و حداقل دو ماشین مجازی کاملاً بیکار باشند، دستور حذف یکی از ماشین های مجازی بیکار صادر می شود. نحوه عملکرد خروجی این الگوریتم

این کمبود به چه صورت است. پس از مشخص شدن حالت سیستم، الگوریتم برای استتباط فازی از پایگاه قاعده‌ای استفاده می‌کند که تصمیم‌گیری‌های آن متناسب با شرایط به وجود آمده در محیط رایانش ابری است. همچنین، الگوریتم در هنگام پیاده‌سازی تصمیم مقیاس‌پذیری، با توجه به این که محیط رایانش ابری در حالت شدید یا حالت عادی باشد (چه در کمبود چه در تامین بیش از حد)، تصمیم به انتخاب قدم مقیاس‌پذیری می‌گیرد. در حالت‌های شدید، الگوریتم برای مقیاس‌پذیری از قدم بزرگتری نسبت به حالت‌های عادی استفاده می‌کند. همان‌طور که در بخش ارزیابی گفته شد، نوآوری معرفی شده در رهیافت و پویایی اضافه شده به روش استتباط فازی، در عملکرد رهیافت بسیار موثر است و تصمیم‌گیری‌های بهینه مناسب و متناسب با محیط رایانش ابری در حال حاضر را عملی می‌سازد. با توجه به این که روش پیشنهادی یک روش مقیاس‌پذیری مبتنی بر استتباط فازی است، با تغییرات در مؤلفه‌ها و ورودی‌های سیستم، می‌توان عملکرد کلی را بهبود داد. به‌طور مثال در رهیافت پیشنهادی برای بررسی و شناخت شرایط محیط از دو مؤلفه MCU و SLA استفاده شده است که برای بهبود عملکرد الگوریتم می‌توان مؤلفه‌های بیشتری مانند مقدار حافظه، پهنای باند شبکه و ویژگی زمانی حال حاضر (به‌طور مثال روزهای آخر هفته یا ساعات اوج درخواست کاربران) را برای بررسی شرایط سیستم در نظر گرفت. همچنین در روش پیشنهاد شده، خروجی سیستم استتباط فازی یک متغیر است که نشان دهنده نحوه مقیاس‌پذیری خواهد بود. برای بهبود بیشتر عملکرد، می‌توان خروجی سیستم را میان دو آستانه بالایی و پایینی برای یک مؤلفه یا چند مؤلفه مورد بررسی، مانند CPU مصرفی، در نظر گرفت و پس از تعیین آستانه‌ها، در این روش مانند روش‌های مبتنی بر آستانه عمل کرد. با تغییر در سیستم استتباط فازی رهیافت پیشنهادی، مانند تغییر توابع عضویت و اصطلاحات فازی، افزایش تعداد متغیرهای ورودی سیستم استتباط فازی، تغییر تعداد پایگاه‌های قاعده و یا افزایش و تغییر قوانین موجود در

پایگاه‌های قواعد موجود نیز عملکرد روش ارائه شده می‌تواند بهبود یابد.

مراجع

- [1] R. Buyya, M. A. Rodriguez, A. N. Toosi, and J. Park, "Cost-efficient orchestration of containers in clouds: A vision, architectural elements, and future directions," 2018, vol. 1108: IOP Publishing, 1 ed., p. 012001.
- [2] V. Persico, D. Grimaldi, A. Pescape, A. Salvi, and S. Santini, "A fuzzy approach based on heterogeneous metrics for scaling out public clouds," IEEE Transactions on Parallel and Distributed Systems, vol. 28, no. 8, pp. 2117-2130, 2017.
- [3] F. Lombardi, A. Muti, L. Aniello, R. Baldoni, S. Bonomi, and L. Querzoni, "PASCAL: An architecture for proactive auto-scaling of distributed services," Future Generation Computer Systems, vol. 98, pp. 342-361, 2019.
- [4] P. Singh, P. Gupta, K. Jyoti, and A. Nayyar, "Research on auto-scaling of web applications in cloud: survey, trends and future directions," Scalable Computing: Practice and Experience, vol. 20, no. 2, pp. 399-432, 2019.
- [5] C. Qu, R. N. Calheiros, and R. Buyya, "A reliable and cost-efficient auto-scaling system for web applications using heterogeneous spot instances," Journal of Network and Computer Applications, vol. 65, pp. 167-180, 2016.
- [6] «Cloud Autoscaling Explained.» <https://www.densify.com/articles/autoscaling>, Last visite: 11/25/2021.
- [7] «Mamdani and Sugeno Fuzzy Inference Systems.» <https://de.mathworks.com/help/fuzzy/types-of-fuzzy-inference-systems.html> (accessed 24 October 2021).
- [8] M. Ghoabaei-Arani, M. Rezaei, and A. Souri, "An auto-scaling mechanism for cloud-based multimedia storage systems: a fuzzy-based elastic controller," Multimedia Tools and Applications, pp. 1-23, 2021.
- [9] C.L. Hung, Y.C. Hu, and K.C. Li, "Auto-scaling model for cloud computing system," International Journal of Hybrid Information Technology, vol. 5, no. 2, pp. 181-186, 2012.
- [10] K. Cetinski and M. B. Juric, "AME-WPC: Advanced model for efficient workload prediction in the cloud," Journal of Network and Computer Applications, vol. 55, pp. 191-201, 2015.
- [11] C.W. Tseng, M.S. Tsai, Y.T. Yang, and L.D. Chou, "A Rapid Auto-Scaling Mechanism in Cloud Computing Environment," Computing, vol.21, 2021.
- [12] Y. Wei, D. Kudenko, S. Liu, L. Pan, L. Wu, and X. Meng, "A reinforcement learning based auto-scaling approach for SaaS providers in dynamic cloud environment," Mathematical Problems in Engineering, vol. 2019, 2019.
- [13] H. Arabnejad, C. Pahl, P. Jamshidi, and G. Estrada, "A comparison of reinforcement learning techniques for fuzzy cloud auto-scaling," in Proceedings of the 17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), Madrid, Spain, 14-17 May 2017, pp. 64-73.
- [14] W.H. Liao, S.C. Kuai, and Y.R. Leau, "Auto-scaling strategy for amazon web services in cloud computing," in Proceedings of the IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity), Chengdu, China, 19-21 December 2015, pp. 1059-1064.