

## تخصیص منابع پردازشی لبه در شبکه‌های اینترنت اشیا

علی سرورامینی

دانشجوی دانشکده مهندسی برق و کامپیوتر، پردیس فنی، دانشگاه تهران، ایران  
پست الکترونیکی: ali.sorouramini@ut.ac.ir

وحید شاه‌منصوری\*

استادیار دانشکده مهندسی برق و کامپیوتر، پردیس فنی، دانشگاه تهران، ایران  
پست الکترونیکی: vmansouri@ut.ac.ir

### چکیده

پردازشی لبه می‌تواند به چند سرویس اختصاص پیدا کند. مسئله به صورت یک مسئله بهینه‌سازی مدل می‌شود که هدف آن بیشینه کردن مجموع سود سرویس‌ها است. مسئله بهینه‌سازی نهایی یک مسئله برنامه‌ریزی غیرخطی عدد صحیح مخلوط است که در حالت کلی به سختی حل می‌شود. الگوریتم زیربهبینه ارائه شده برای حل این مسئله جوابی قابل قبول برای آن به دست می‌دهد که به صورت توزیع شده قابل اجرا است.

**واژه‌های کلیدی:** اینترنت اشیا، رایانش لبه، تخصیص

منابع

مقدمه

گسترش اینترنت اشیا و موفقیت سرویس‌های ابری باعث ایجاد الگوی جدیدی در پردازش داده‌ها به نام رایانش لبه<sup>۱</sup> شده است. طبق برآوردهای انجام شده در [۱] تعداد دستگاه‌های متصل شده به شبکه ۵۰ میلیارد عدد خواهد بود. بعضی از کاربردهای اینترنت اشیا نیاز به زمان پاسخ کوتاه دارند، بعضی ممکن است دارای داده‌های

با حرکت به سمت عصر اینترنت اشیا، تعداد دستگاه‌های متصل به اینترنت به صورت نمایی در حال افزایش است. تعداد زیاد دستگاه‌های متصل باعث ایجاد گلوگاه‌هایی در حوزه‌های مختلف مانند اتصال دستگاه‌ها، انتقال و پردازش داده‌ها می‌شود. پردازش لبه یک روش مناسب برای پردازش حجم زیاد داده‌های تولید شده توسط اشیا متصل به اینترنت به شمار می‌رود. به لطف پیشرفت‌های ایجاد شده در فناوری‌های مجازی‌سازی، دروازه‌های شبکه و دستگاه‌های لبه شبکه می‌توانند ظرفیت پردازشی اضافی خود را در اختیار سرویس‌های اینترنت اشیا قرار دهند. تعداد بسیار زیاد دستگاه‌های لبه و سرویس‌های شبکه اینترنت اشیا، باعث پیچیده شدن مسئله تخصیص منابع پردازشی به سرویس‌های شبکه اینترنت اشیا می‌شود. در این مقاله، وظیفه تخصیص منابع پردازشی مورد نیاز سرویس‌های اینترنت اشیا را برای تعداد زیاد سرویس‌ها در نظر می‌گیریم. هر سرویس می‌تواند از چند منبع پردازشی استفاده کند و هر منبع

\* نویسنده مسئول

1- Edge Computing

محرمانه و شخصی باشند و بعضی از این کاربردها می‌توانند بار سنگینی برای شبکه داشته باشند و پردازش ابری ممکن است روش مناسبی برای این کاربردها نباشد. در اینجا منظور از لبه، همه منابع پردازشی و شبکه‌ای است که بین منبع داده‌ها (جایی که داده‌ها در آن جا تولید می‌شوند) و مراکز داده ابری قرار دارند. برای مثال یک دروازه شبکه در یک خانه هوشمند، می‌تواند یک منبع پردازشی لبه بین اشیاء و مرکز داده ابری باشد یا یک مرکز داده کوچک، می‌تواند یک منبع پردازشی لبه بین دستگاه‌های سیار و ابر در نظر گرفته شود. منطق در رایانش لبه این است که پردازش داده‌ها باید در همسایگی منبع داده‌ها انجام شود. با این منطق رایانش لبه و رایانش مه دو مفهوم یکسان را خواهند داشت با این تفاوت که رایانش لبه تمرکزش سمت اشیاء است ولی رایانش مه تمرکزش سمت زیرساخت است. فشار از سمت رایانش ابری، تغییر از مصرف کننده داده به تولید کننده داده و کشش از سمت اینترنت اشیاء از عوامل جذابیت رایانش لبه است [۲].

در الگوی رایانش لبه، اشیاء نه تنها مصرف کننده داده بلکه تولید کننده داده نیز هستند. در واقع در لبه، اشیاء نه تنها می‌توانند سرویس و محتوا از ابر درخواست نمایند بلکه می‌توانند وظایف پردازشی را هم انجام دهند. در [۳] برنامه‌ای برای تشخیص چهره ساخته شده است که با انتقال پردازش از ابر به لبه شبکه، زمان پاسخ برنامه از ۹۰۰ به ۱۶۰ میلی ثانیه کاهش پیدا کرده است. نویسندگان در [۴] به بررسی یک بُن‌سازه رایانش مه که از Docker استفاده می‌کند پرداخته‌اند. بررسی آن‌ها نشان می‌دهد که این بُن‌سازه کارایی بهتری نسبت به روش‌های فعلی دارد. در [۵] نویسندگان از واحدهای پردازشی ابری کوچک<sup>۲</sup> برای پردازش وظایف دستیار شناختی پوشیدنی استفاده کرده‌اند و نتایج نشان می‌دهد که زمان پاسخ بین ۸۰ تا ۲۰۰ میلی ثانیه کاهش پیدا کرده است. در [۶] نویسندگان، Docker را به عنوان فناوری که امکان استفاده از پردازش لبه را

2- Cloudlet

فراهم می‌کند بررسی کرده‌اند. در [۷] و [۸] از فناوری‌های مجازی‌سازی سبک برای طراحی دروازه‌های اینترنت اشیاء استفاده شده است. در [۷] از یک نرم‌افزار مجازی‌سازی برای استقرار انبوه سرویس‌ها در لبه شبکه استفاده کرده‌اند. نکته قابل توجه در تحلیل آن‌ها، امکان اثر متقابل بین حسگرهای اینترنت اشیاء و دروازه‌ها است. تحلیل آن‌ها نشان می‌دهد که تخصیص پویای سرویس‌ها با استفاده از کانتینرها مزایایی از بعد کارایی در دروازه‌ها دارد.

[۹] به بررسی کارآیی فناوری‌های کانتینری روی کامپیوترهای تک‌بوردی مختلف می‌پردازد. در همه این مقالات نتیجه‌گیری بر این است که با استفاده از کانتینرها در دستگاه‌های مورد نظر، سربار پردازشی قابل چشم پوشی است. تخصیص منابع پردازشی برای سرویس‌ها در اینترنت اشیاء به دلیل تعداد زیاد سرویس‌ها و منابع پردازشی، یک مسئله پیچیده است. نویسندگان در [۱۰] چهارچوبی پیاده‌سازی کرده‌اند که در اختصاص منابع پردازشی در یک شبکه ابری و مه، تأخیر انتقال و مصرف انرژی را بهینه می‌کند. در [۱۱] نویسندگان مسئله توزیع پردازش سرویس‌ها را به صورت مسئله بهینه‌سازی برای کمینه کردن هزینه فرمول‌بندی کرده‌اند. [۱۲] به بررسی تخصیص منابع پردازشی در یک شبکه رایانش مه می‌پردازد. در این مقاله شبکه با سه لایه مدل می‌شود که شامل لایه گره‌های مه، لایه اپراتورهای مراکز داده و لایه کاربران است و از بازی استکلبرگ<sup>۳</sup> و نظریه تطبیق<sup>۴</sup> برای پیدا کردن نقطه تعادل استفاده شده است. در [۱۳] از رایانش مه برای کاهش بار پردازشی مراکز داده استفاده شده است. مسئله به صورت یک بهینه‌سازی فرمول‌بندی شده است که تلاش می‌کند تأخیر، مصرف انرژی و هزینه را کمینه کند. نویسندگان در [۱۴] مسئله تخصیص منابع رادیویی و پردازشی را بررسی کرده‌اند. در این مقاله تابع هدف شامل بیشینه کردن نرخ کاربران و کمینه کردن هزینه است و از نظریه تطبیق برای پیدا کردن جواب استفاده شده است.

3- Wearable Cognitive Assistance

4- Gateway

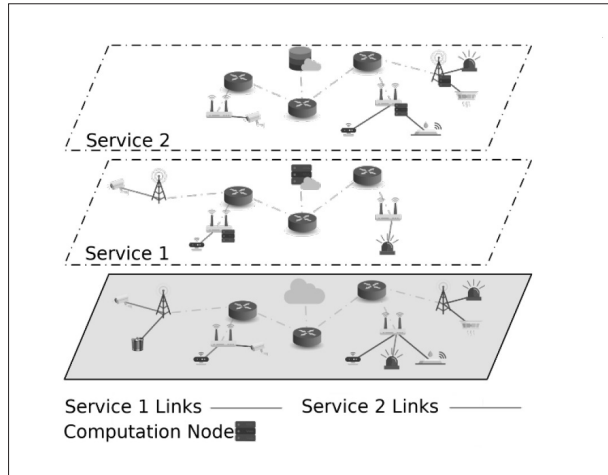
5- Stackelberg Game

منابع پردازشی قابل استفاده برای سرویس  $S$  با مجموعه  $C_S$  نشان داده می شود که زیرمجموعه ای از  $C$  است. فراهم کننده های زیرساخت<sup>۷</sup> علاقه دارند که توان پردازشی خود را با دیگران به اشتراک بگذارند. مدل سیستم در شکل ۱ نشان داده شده است.

$\phi_i$  نشان دهنده توان پردازشی فراهم کننده زیرساخت  $i$  است و  $v_i$  درصدی از توان پردازشی که فراهم کننده زیرساخت پردازشی  $i$  مایل به اشتراک گذاری است را نشان می دهد. در نتیجه، فراهم کننده زیرساخت پردازشی  $i$  مایل به اشتراک گذاری توان پردازشی برابر با  $\phi_i v_i$  خواهد بود. در اینجا، گره های لبه<sup>۸</sup>، گره هایی هستند که زیرساخت پردازشی را در لبه شبکه فراهم می کنند. در مقابل، فراهم کننده های زیرساخت ابری، منابع پردازشی را در فاصله ای دورتر از سرویس ها فراهم می کنند. بنابراین، گره های لبه، توان پردازشی کمتر با تأخیر کمتر را در مقایسه با فراهم کننده های زیرساخت ابری مهیا می کنند. بدهستان، بین میزان توان پردازشی و تأخیر در پردازش، باعث پیچیده شدن انتخاب گره های پردازشی برای سرویس ها می شود.

$\mu_{i,s}$  نشان دهنده نرخ سرویس منبع پردازشی  $i$  برای سرویس  $S$  است.  $r_s$  نرخ انتخابی ارسال نمونه های سرویس  $S$  است و  $R_s$  نرخ مطلوب ارسال نمونه ها برای سرویس  $S$  را نشان می دهد. جدول ۱ به صورت خلاصه پارامترهای استفاده شده در این مقاله را معرفی می کند.

در این جا فرض می کنیم که ارزش منبع پردازشی  $i$  برای سرویس  $S$  از دو بخش تشکیل شده است. بخش اول، تابع نرخ ارسال نمونه ها توسط حسگرهای سرویس  $S$  به منبع پردازشی  $i$  است. این تابع را  $f_r$  نام گذاری می کنیم و فرض می کنیم که  $f_r$  تابعی از نرخ ارسال نمونه ها توسط حسگرهای سرویس  $S$  و نرخ مطلوب آن سرویس است. در واقع  $f_r$  بیان کننده تأثیر  $r_s$  در ارزش به دست آمده از پردازش اطلاعات برای سرویس ها است. بخش دوم هم



شکل ۱: مدل سیستم

در این مقاله، مسئله تخصیص منابع برای سرویس های اینترنت اشیا در نظر گرفته می شود. فرض بر این است که با کمک فناوری Docker هر منبع پردازشی لبه می تواند به چند سرویس اینترنت اشیا سرویس دهد. همچنین هر سرویس می تواند از چند منبع پردازشی استفاده کند. مسئله تخصیصی منابع به صورت یک مسئله بهینه سازی مدل می شود. سپس یک روش زیربینه برای حل مسئله با پیچیدگی پردازش پایین ارائه می شود.

ساختار مقاله در ادامه به صورت زیر است. در فصل دوم به بررسی مدل سیستم می پردازیم و مسئله بهینه سازی تخصیص منابع را در قالب یک مسئله بهینه سازی معرفی می کنیم. فصل سوم روش زیربینه پیش نهادی را معرفی می کند. فصل چهارم نتایج شبیه سازی را نشان می دهد و فصل پنجم نتیجه گیری مقاله را شامل می شود.

## ۲- مدل سیستم

در یک منطقه جغرافیایی تعداد  $N$  سرویس اینترنت اشیا وجود دارد و تعدادی فراهم کننده زیرساخت<sup>۶</sup> که منابعی پردازشی را در اختیار سرویس ها قرار می دهند. هر سرویس از تعدادی حسگر و عملگر تشکیل شده که نیاز به پردازش اطلاعات دارند. مجموعه سرویس ها را با  $S$  و مجموع منابع پردازشی را با  $C$  نمایش می دهیم.

7-Infrastructure Provider  
8- Edge Nodes

6- Matching Theory

جدول ۱: نمادهای استفاده شده

نشانه	توضیح
$N$	تعداد سرویس‌های شبکه
$S$	مجموعه سرویس‌های شبکه
$C$	مجموعه منابع پردازشی شبکه
$C_s$	مجموعه منابع پردازشی که سرویس $S$ می‌تواند از آن‌ها استفاده کند.
$\alpha_s$	ضریب اهمیت سرویس $S$
$\mu_{i,s}$	نرخ سرویس منبع پردازشی $i$ برای سرویس $S$
$r_s$	نرخ انتخابی ارسال نمونه‌های سرویس $S$
$r_{i,s}$	نرخ که سرویس $S$ برای پردازش به منبع پردازشی $i$ می‌فرستد
$R_s$	نرخ مطلوب ارسال نمونه‌ها برای سرویس $S$
$\omega_s$	ضریب وزن دهی برای مشخص کردن نسبت اهمیت تأخیر نسبت به اختلاف نرخ انتخابی و نرخ مطلوب برای سرویس‌ها
$\delta_{i,s}$	متغیر دودویی که مشخص می‌کند سرویس $S$ از منبع پردازشی $i$ استفاده می‌کند یا نه
$d_{i,s}$	تأخیر منبع پردازشی $i$ برای سرویس $S$
$d_{i,s}^{CPU}$	تأخیر پردازشی منبع $i$ برای سرویس $S$
$d_{i,s}^{net}$	تأخیر شبکه سرویس $S$ زمانی که از منبع پردازشی $i$ استفاده می‌کند
$d_s^{max}$	بیشترین تأخیر سرویس $S$ بین منابع پردازشی که از آن‌ها استفاده می‌کند
$u_{i,s}$	بخشی از ظرفیت پردازشی منبع پردازشی $i$ که توسط سرویس $S$ استفاده می‌شود
$Q_s$	حداکثر تعداد منابع پردازشی که سرویس $S$ مجاز به استفاده از آن‌ها است
$Q_i$	حداکثر تعداد سرویس‌هایی که منبع پردازشی $i$ می‌تواند به آن‌ها اختصاص پیدا کند
$f_r, f_d$	توابع مطلوبیت تأخیر و سرعت

تابعی از تأخیر ایجاد شده (مجموع تأخیر شبکه و تأخیر پردازش) است که آن را با  $f_d$  نمایش می‌دهیم.  $d_s^{max}$  نشان‌دهنده تأخیر سرویس  $S$  است که در ادامه معرفی می‌کنیم. اگر فرض کنیم  $R_s$  نرخ ارسال مطلوب نمونه‌ها برای سرویس  $S$  و  $r_s$  نرخ ارسال نمونه‌های سرویس  $S$  باشد، ارزشی که سرویس  $S$  با استفاده از منبع پردازشی  $i$  به دست می‌آورد را به صورت زیر می‌توان نوشت

اختصاص یافته برای سرویس اهمیتی ندارد و فقط تأخیر برای سرویس مهم است. واضح است که اگر  $\omega = 0$  باشد، ارزشی که ارزشی به دست آمده برای سرویس را بیشینه می‌کند  $r_s = 0$  است. در نتیجه  $\omega = 0$  نتیجه قابل قبولی خواهد داشت.  $\alpha_s$  برای اثر دادن ارزش ذاتی سرویس‌ها است. پس باید برای سرویس با ارزش بالاتر،  $\alpha_s$  مقدار بیشتری داشته باشد. در ادامه برای نشان دادن اختصاص منبع پردازشی  $i$  به سرویس  $S$  از  $\delta_{i,s}$  استفاده می‌کنیم.  $\delta_{i,s}$  برابر یک است اگر منبع پردازشی  $i$  به سرویس  $S$  اختصاص پیدا کرده باشد. در غیر این صورت  $\delta_{i,s}$  برابر صفر است.

$$U_s = \omega_s f_r(r_s, R_s) + (1 - \omega_s) f_d(d_s^{max}). \quad (1)$$

نکته قابل توجه در این رابطه این است که  $\omega \in [0, 1]$  یک ضریب وزن دهنده است که تعیین کننده نسبت تاثیر  $f_r$  به  $f_d$  است. اگر  $\omega = 1$  باشد ضریب  $f_d$  صفر خواهد بود، در نتیجه تأخیر اثری در ارزش به دست آمده برای سرویس نخواهد داشت. اما  $\omega = 0$  بیانگر این است که نرخ

علوم رایانشی / تابستان ۱۴۰۰

منبع پردازشی  $i$  برای سرویس  $S$  پردازش می‌شود در نظر بگیریم، رابطه زیر برای  $r_s$  برقرار است

$$r_s = \sum_{i \in C_s} r_{i,s}. \quad (2)$$

فرض بر این است که هر سرویس می‌تواند از چند منبع پردازشی استفاده کند، بیشترین تأخیر منابع پردازشی مورد استفاده هر سرویس را به عنوان تأخیر آن سرویس در نظر می‌گیریم. اگر فرض کنیم  $d_{i,s}$  تأخیر سرویس  $S$  در منبع پردازشی  $i$  و  $d_s^{max}$  بیشترین تأخیر سرویس  $S$  باشد،  $d_s^{max}$  باید از همه  $d_{i,s}$ ها بزرگتر باشد. در نتیجه قید زیر باید برای همه منابع پردازشی انتخاب شده توسط سرویس  $S$  برقرار باشد

$$d_s^{max} = \max_i \{d_{i,s}\} \geq d_{i,s}; \forall s \in S, i \in C_s, \delta_{i,s} = 1. \quad (3)$$

هر کدام از تأخیرهای  $d_{i,s}$ ها از دو بخش تشکیل شده‌اند. قسمت اول آن مربوط به تأخیر شبکه است که برابر است با زمانی که طول می‌کشد تا نمونه‌ها به منبع پردازشی برسند. به علاوه زمانی که طول می‌کشد نتیجه به مقصد برسد. برای سرویس  $S$  تأخیر شبکه زمانی که از منبع  $i$  استفاده می‌شود با  $d_{i,s}^{net}$  نشان داده می‌شود. قسمت دوم، تأخیر پردازش نمونه‌ها در منابع پردازشی است. این تأخیر برابر زمانی است که طول می‌کشد تا نمونه‌ها پس از رسیدن به منابع پردازشی، پردازششان پایان یابد. برای محاسبه تأخیر پردازشی از نظریه صف استفاده می‌کنیم. به دلیل این‌که فرض بر این است که منابع پردازشی بین چند سرویس ممکن است تقسیم شوند، از مدل  $M/M/1$  استفاده می‌کنیم چرا که در این حالت پردازنده منابع پردازشی بین سرویس‌های آن منبع به اشتراک گذاشته می‌شود. در مدل  $M/M/1$  میانگین تأخیر برابر است با  $1/(\mu - \lambda)$  [15].

ظرفیت پردازشی منبع پردازشی  $i$  را با  $\varphi_i = \phi_i v_i$  نشان می‌دهیم. متغیر  $u_{i,s}$  تعیین می‌کند که چه درصد از ظرفیت پردازشی منبع پردازشی  $i$  به سرویس  $S$  اختصاص پیدا می‌کند. واضح است که سرویس‌ها نمی‌توانند بیش‌تر از ظرفیت منابع پردازشی از آن‌ها استفاده کنند. در نتیجه رابطه زیر برای مقادیر ظرفیت‌های پردازشی تخصیص

یافته به سرویس‌ها برای همه منابع پردازشی باید برقرار باشد

$$\sum_{s \in S_i} u_{i,s} \leq 1, \forall i \in C. \quad (4)$$

این رابطه بیان می‌کند که مجموع ظرفیت پردازشی اختصاص داده شده به سرویس‌ها برای هر منبع پردازشی نباید از ظرفیت کلی آن منبع پردازشی بیشتر باشد. نرخ سرویس برای سرویس  $S$  در منبع پردازشی  $i$  از رابطه زیر به دست خواهد آمد

$$\mu_{i,s} = \frac{\varphi_i}{F_s} u_{i,s} = \zeta_{i,s} u_{i,s}, \forall s \in S, i \in C_s. \quad (5)$$

در واقع  $\mu_{i,s}$ ، نرخ اجرا شدن  $F_s$  دستورالعمل پردازنده توسط قسمتی از کل ظرفیت پردازنده منبع پردازشی  $i$  که به وسیله  $u_{i,s}$  مشخص می‌شود است. در این رابطه  $\zeta_{i,s}$  یک عدد ثابت است که مقدار آن برابر است با  $\frac{\varphi_i}{F_s} u_{i,s}$ . با این تفاسیر رابطه زیر را برای تأخیر پردازشی سرویس  $S$  در منبع پردازشی  $i$  وقتی از آن منبع پردازشی استفاده می‌کند می‌توان نوشت

$$d_{i,s}^{CPU} = \frac{1}{\zeta_{i,s} u_{i,s} - r_{i,s}}. \quad (6)$$

در نتیجه رابطه زیر برای تأخیر سرویس  $S$  در منبع پردازشی  $i$  برقرار است

$$d_{i,s} = d_{i,s}^{net} + \frac{1}{\zeta_{i,s} u_{i,s} - r_{i,s}}. \quad (7)$$

این رابطه به عنوان قید در یک بهینه‌سازی محدب قابل استفاده نیست. به همین دلیل سعی می‌کنیم تغییراتی در آن ایجاد کنیم تا به یک قید محدب تبدیل شود. با توجه به وجود  $d_{i,s}^{max}$  در تابع هدف بهینه‌سازی، شرایط تابع  $f_d$  و این‌که رابطه (۳) جزء قیدهایی مسئله است، می‌توانیم علامت = را با علامت  $\leq$  جایگزین کنیم.

$$\frac{1}{\zeta_{i,s} u_{i,s} - r_{i,s}} \leq d_{i,s} - d_{i,s}^{net}. \quad (8)$$

چون لگاریتم یک تابع صعودی است، می‌توان از دو طرف نامساوی لگاریتم گرفت و جهت علامت نامساوی تغییری نمی‌کند. با لگاریتم گرفتن از طرفین می‌توان قید مناسب برای استفاده در بهینه‌سازی محدب را به دست آورد

$$-\log(\zeta_{i,s} u_{i,s} - r_{i,s}) - \log(d_{i,s} - d_{i,s}^{net}) \leq 0. \quad (9)$$

دلیل محدب بودن رابطه (۹) این است که  $\zeta_{i,s} u_{i,s} - r_{i,s}$  و  $d_{i,s} - d_{i,s}^{net}$  هر دو تابع محدب هستند. همچنین تابع  $-\log$  هم یک تابع محدب است. با توجه به قضیه ترکیب توابع محدب [۱۶]، (۹) یک قید محدب است.

تابع مطلوبیت سرویس  $\mathbf{s}$ ، ترکیبی وزن دار از تأخیر و نرخ پردازش اختصاص داده شده به سرویس است که در فرمول رابطه (۱۰) ارائه شده است. رابطه (۱۰) تابع هدف بهینه‌سازی را نشان می‌دهد که مجموع وزن دار برای سرویس‌های مختلف می‌تواند متفاوت باشد و اهمیت یک سرویس را در مقایسه بقیه سرویس‌ها نشان دهد.

$$\sum_{s=1}^N \alpha_s (\omega_s f_r(r_s, R_s) + (1 - \omega_s) f_d(d_s^{max})). \quad (10)$$

با توجه به آنچه که تا این جا گفته شد، می‌توان مسئله

بهینه‌سازی را به صورت زیر نوشت

$$\begin{aligned} & \max_{r_{i,s}, u_{i,s}, \delta_{i,s}} \alpha_s (\omega_s f_r(r_s, R_s) + (1 - \omega_s) f_d(d_s^{max})) \\ & \text{Subject to} \\ & r_s = \sum_{i \in S_i} r_{i,s}, \forall s \in S \quad (11a) \\ & r_{i,s} \leq \eta \zeta_{i,s} u_{i,s}, \forall s \in S, i \in C_s \quad (11b) \\ & u_{i,s} \leq \delta_{i,s}, \forall s \in S, i \in C_s \quad (11c) \\ & d_{i,s} \leq d_s^{max}, \forall s \in S, i \in C_s, \text{ if } \delta_{i,s} = 1 \quad (11d) \\ & -\log((\zeta_{i,s} u_{i,s} - r_{i,s})(d_{i,s} - d_{i,s}^{net})) \leq 0, \forall s \in S, i \in C_s, \text{ if } \delta_{i,s} = 1 \quad (11e) \\ & \sum_{s \in S_i} u_{i,s} \leq 1, \forall i \in C \quad (11f) \\ & \sum_{s \in S_i} \delta_{i,s} \leq Q_i, \forall i \in C \quad (11g) \\ & \sum_{i \in C_s} \delta_{i,s} \leq Q_s, \forall s \in S \quad (11h) \\ & r_{i,s} \geq 0, \forall s \in S, i \in C_s \quad (11i) \\ & 0 \leq u_{i,s} \leq 1, \forall s \in S, i \in C_s \quad (11j) \\ & \delta_{i,s} \in \{0,1\}, \forall s \in S, i \in C_s \quad (11k) \end{aligned}$$

در این مسئله قید (۱۱ج) باعث می‌شود که زمانی که سرویس  $s$  از منبع پردازشی  $i$  استفاده نمی‌کند  $u_{i,s} = 0$  باشد و در صورت استفاده  $1 \leq u_{i,s}$  باشد. باید توجه کرد که قیدهای (۱۱ط) و (۱۱ب) باعث می‌شوند که  $u_{i,s}$ ‌ها مثبت باشند. قید (۱۱و) برای این در بهینه‌سازی حضور دارد که میزان استفاده از منابع پردازشی نمی‌تواند بیشتر از ظرفیت آن‌ها باشد. قید (۱۱ز) نشان‌دهنده حداکثر تعداد سرویس‌هایی است که یک منبع پردازشی می‌تواند به آن‌ها اختصاص پیدا کند. قید (۱۱ح) برای محدود کردن حداکثر تعداد منابع پردازشی که یک سرویس می‌تواند استفاده کند است.

این مسئله بهینه‌سازی هم یک مسئله برنامه‌ریزی غیرخطی عدد صحیح مخلوط است که پیدا کردن جواب بهینه آن ساده نیست. به همین دلیل در بخش بعدی یک الگوریتم برای پیدا کردن جواب زیر بهینه آن معرفی می‌کنیم.

### ۳- الگوریتم زیر بهینه

از آنجایی که پیدا کردن مقدار بهینه مسئله (۱۱) نیاز به زمان طولانی دارد، برای پیدا کردن جواب آن یک الگوریتم مبتنی بر تکرار معرفی می‌کنیم که در زمان کوتاهی به جواب می‌رسد. در هر تکرار الگوریتم، برای هر سرویس، یک تغییر در منابع پردازشی اختصاص یافته به آن سرویس ایجاد می‌کنیم. این تغییر می‌تواند حذف کردن، اضافه کردن یا جابه‌جایی منابع پردازشی آن سرویس باشد به شرطی که پس از تغییر، قیدهای حداکثر تعداد سرویس‌های منابع پردازشی و حداکثر تعداد منابع پردازشی سرویس برقرار باشند. این کار را در هر تکرار برای همه سرویس‌ها انجام می‌دهیم و تغییر ایجاد شده در تابع هدف بهینه‌سازی با مرحله قبل را در نظر می‌گیریم. پس از انجام این کار، تغییری که بیشترین افزایش را در تابع هدف بهینه‌سازی دارد به تخصیص منابع اعمال می‌کنیم. این الگوریتم تا زمانی که مقدار افزایش تابع هدف بهینه‌سازی در هر مرحله بیشتر از مقدار مشخص  $\epsilon$  باشد ادامه می‌دهیم.

الگوریتم ۱ به صورت خلاصه مراحل این الگوریتم را نشان می‌دهد. خط ۳ باعث می‌شود که الگوریتم تا زمانی که تغییرات مرحله قبل بیش از  $\epsilon$  باشد، ادامه پیدا کند. سپس در خط ۴ مقدار تغییرات در مرحله جدید برابر صفر قرار داده می‌شود. در خط ۵ همه سرویس‌ها مورد بررسی قرار می‌گیرند. برای بررسی هر سرویس در خط ۶، همه منابع پردازشی ممکن  $C_1$  و  $C_2$  بررسی می‌شوند.  $C_1$  منبع پردازشی است که از مراحل قبل به سرویس مورد بررسی اختصاص پیدا کرده و حذف شدن آن را بررسی می‌کنیم و  $C_2$  منبع پردازشی است که اختصاص پیدا کردن



الگوریتم ۱: الگوریتم تخصیص منابع پردازشی مبتنی بر مزایده

```

1:  $\nu \leftarrow -\infty$ 
2:  $d \leftarrow \infty$ 
3: while  $d > \epsilon$  do
4:    $d \leftarrow 0$ 
5:   for  $s \in S$  do
6:     for  $(c_1, c_2) \in (I_s \cup \{-\}) \times ((C_s \setminus I_s) \cup \{-\})$ 
7:       do
8:         if  $s \in S_{c_2}$  then
9:            $I_s \leftarrow (I_s \setminus \{c_1\}) \cup \{c_2\}$ 
10:           $I_{c_1} \leftarrow I_{c_1} \setminus \{s\}$ 
11:           $I_{c_2} \leftarrow I_{c_2} \cup \{s\}$ 
12:          if  $|I_s| \leq Q_s$  and  $|I_{c_2}| \leq Q_{c_2}$  then
13:             $\nu' = \text{Optimal Value}$ 
14:            if  $\nu' > \nu$  then
15:               $d \leftarrow d + \nu' - \nu$ 
16:               $\nu \leftarrow \nu'$ 
17:               $s^{\text{new}} \leftarrow s$ 
18:               $c_1^{\text{new}} \leftarrow c_1$ 
19:               $c_2^{\text{new}} \leftarrow c_2$ 
20:            end if
21:          end if
22:           $I_s \leftarrow (I_s \setminus \{c_2\}) \cup \{c_1\}$ 
23:           $I_{c_2} \leftarrow I_{c_2} \setminus \{s\}$ 
24:           $I_{c_1} \leftarrow I_{c_1} \cup \{s\}$ 
25:        end if
26:      end for
27:    end for
28:     $I_s^{\text{new}} \leftarrow (I_s^{\text{new}} \setminus \{c_1^{\text{new}}\}) \cup \{c_2^{\text{new}}\}$ 
29:     $I_{c_1}^{\text{new}} \leftarrow I_{c_1}^{\text{new}} \setminus \{s^{\text{new}}\}$ 
30:     $I_{c_2}^{\text{new}} \leftarrow I_{c_2}^{\text{new}} \cup \{s^{\text{new}}\}$ 
31:  end while

```

به روز رسانی می شوند (خطوط ۱۴ تا ۱۸). در پایان هر مرحله تخصیص منابعی که بیشتری افزایش را در تابع هدف مسئله داشته، اعمال می شود (خطوط ۲۷ تا ۲۹).

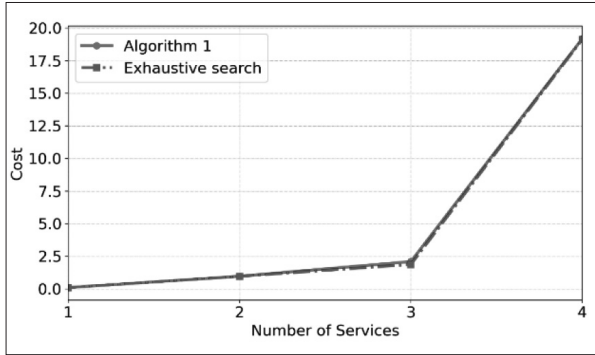
در ادامه به بررسی پیچیدگی الگوریتم زیربینه می پردازیم. در ابتدای شروع الگوریتم، نرخ انتخابی همه سرویس ها صفر است و هیچ منبع پردازشی اختصاص پیدا نکرده است. در نتیجه مقدار تابع هدف بهینه سازی  $-\sum_{s=1}^N \alpha \omega_s R$  است. علاوه بر این تابع هدف بهینه سازی همواره یک عدد منفی است. با در نظر گرفتن این که در هر تکرار، مقدار تابع هدف حداقل به اندازه  $\epsilon$  افزایش پیدا می کند، تعداد تکرارهای الگوریتم نمی تواند بیش از  $\sum_{s=1}^N \alpha \omega_s R / \epsilon$  باشد.

اگر فرض کنیم در مرحله ای دلخواه از این الگوریتم سرویس  $S$  به  $m$  منبع پردازشی اختصاص پیدا کرده باشد و  $m'$  تعداد منابع پردازشی اختصاص پیدا نکرده به سرویس  $S$  باشد، باید  $(m+1) \times (m'+1) - 1$  بار مسئله بهینه سازی برای سرویس  $S$  حل شود. بنابراین، می توان نتیجه گرفت که تعداد دفعاتی که مسئله بهینه سازی در هر بار تکرار الگوریتم حل می شود، کمتر از  $N(M+1)^2$  است. در نتیجه تعداد دفعاتی که مسئله بهینه سازی حل می شود کمتر از  $(M+1)^2 \sum_{s=1}^N \alpha_s \omega_s R_s / \epsilon$  است. از این رابطه واضح است که با افزایش  $\epsilon$  انتظار می رود الگوریتم زودتر پایان یابد. باید توجه داشت که  $(m+1) \times (m'+1) - 1$  حد بالای تعداد دفعات حل شدن مسئله بهینه سازی برای هر سرویس را نشان می دهد و محدود بودن انتخاب های سرویس ها و منابع پردازشی (مجموعه های  $S$  و  $C$  برای هر سرویس و منبع پردازشی) باعث می شود که تعداد واقعی، کمتر از این تعداد باشد.

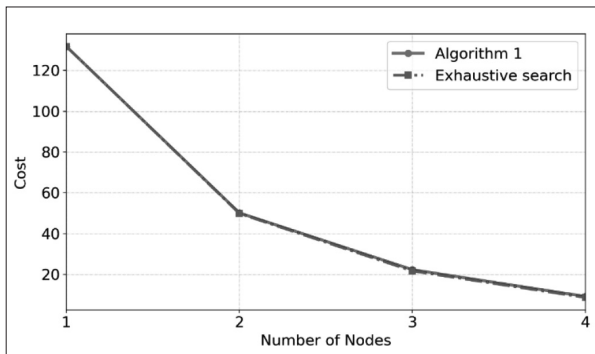
#### ۴- نتایج شبیه سازی

یک شبکه با تعداد  $N$  سرویس و  $M$  منبع پردازشی که در داخل یک دایره به شعاع  $R = 100m$  به صورت تصادفی پراکنده شده اند را در نظر می گیریم. برای محاسبه تأخیر،

آن به سرویس را بررسی می کنیم که هرکدام می تواند تهی باشد. در واقع در هر بار بررسی، منبع پردازشی  $C_1$  را از مجموعه منابع پردازشی اختصاص پیدا کرده به سرویس حذف کرده و  $C_2$  را به آن اضافه می کنیم. اگر منبع پردازشی  $C_2$  تهی باشد یعنی فقط منبع پردازشی  $C_1$  از منابع اختصاص پیدا کرده به سرویس حذف می شود و اگر  $C_1$  تهی باشد یعنی فقط منبع پردازشی  $C_2$  به سرویس اختصاص پیدا می کند. برای این کار در خطوط ۸ تا ۱۰ این تغییر را به صورت موقت در اختصاص منابع اعمال می کنیم و در خطوط ۲۱ تا ۲۳ اختصاص منابع را به حالت قبل باز می گردانیم. پس از بررسی معتبر بودن این تغییرات در خط ۱۱، در خط ۱۲ مقدار بهینه تابع هدف با منابع اختصاص یافته جدید محاسبه می شود و نتیجه در  $\nu'$  قرار می گیرد. در خط ۱۳، با  $\nu'$  بیشترین مقداری که تا الان به دست آمده،  $\nu$  مقایسه می شود و در صورتی که  $\nu'$  از آن بیشتر بود، مقادیر  $\nu$  و  $d$  و بهترین تخصیص منابع



شکل ۲: مقایسه مقدار بهینه حاصل از روش جستجوی فراگیر و الگوریتم ۱ وقتی تعداد سرویس‌ها افزایش پیدا می‌کند.



شکل ۳: مقایسه مقدار بهینه حاصل از روش جستجوی فراگیر و الگوریتم ۱ وقتی تعداد گره‌های پردازشی افزایش پیدا می‌کند.

همانطور که در شکل هم مشخص است با افزایش تعداد منابع پردازشی، میانگین هزینه سرویس‌ها کاهش می‌یابد. علاوه بر این، میانگین هزینه سرویس‌های فرآیندهای منبع پردازشی ابری در آن وجود داشت، همواره کمتر از فرآیندهای بدون منبع رایانش ابری است. دلیلش آن است که هرچه تعداد منابع پردازشی بیشتر باشد، ظرفیت پردازشی کل بیشتر می‌شود و برای سرویس‌ها انتخاب بیش‌تری وجود دارد. همچنین، با افزایش تعداد منابع پردازشی لبه شبکه، اختلاف این دو فرآیندها کاهش پیدا می‌کند تا در تعداد ۲۰ منبع پردازشی لبه شبکه، تقریباً صفر می‌شود. دلیلش آن است که با افزایش منابع پردازشی لبه شبکه همه پردازش‌ها در لبه شبکه انجام شده و تاثیر وجود منبع پردازشی ابری در میانگین هزینه‌ها به صفر می‌رسد.

اثر پارامتر  $\omega_s$  در این روش تخصیص منابع در شکل ۵ بررسی شده است. در این شکل اثر  $\omega_s$  روی اختلاف نرخ بهینه و مطلوب و بیش‌ترین تأخیر برای یک سرویس در بازه  $\omega_s \in [30, 70]$  رسم شده است. همان‌طور که بیان شد  $\omega_s$

تأخیر ناشی از مسیریاب‌های شبکه را در نظر می‌گیریم و فرض می‌کنیم این تأخیر بین ۱۵ تا ۲۰ میلی ثانیه است. فرض می‌کنیم ظرفیت پردازشی منابع پردازشی لبه شبکه به صورت تصادفی در محدوده  $[1200, 1400]$  توزیع شده‌اند. برای ظرفیت پردازشی منابع پردازشی ابری هم مقدار ثابت 10000 را در نظر می‌گیریم. برای  $\alpha_s$  فرض می‌کنیم که به صورت تصادفی در بازه  $[1, 10]$  توزیع شده است. برای رسیدن به نتایج معنی دار،  $\beta = (1 - \omega_s) / \omega_s$  را به صورت یکنواخت در بازه  $[30, 70]$  در نظر می‌گیریم. برای شبیه‌سازی توابع زیر را فرض می‌کنیم. قابل ذکر است الگوریتم با تغییر توابع نیز همچنان کارکرد درستی دارد.

$$f_r(r_s, R_s) = -|R_s - r_s| \quad (12)$$

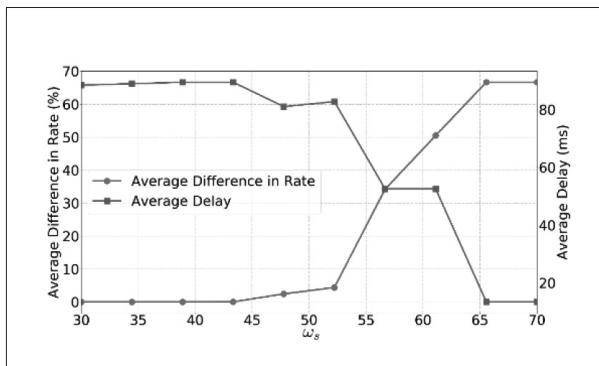
$$f_d(d_{i,s}^{CPU} + d_{i,s}^{net}) = -d_{i,s}^{CPU} - d_{i,s}^{net} \quad (13)$$

برای بررسی کارآمدی الگوریتم مقدار تابع هدف را با روش جستجوی فراگیر برای ۴ سرویس و ۴ گره پردازشی مقایسه می‌کنیم. در شکل ۲ تعداد سرویس‌ها را از ۱ تا ۴ افزایش می‌دهیم و در شکل ۳ تعداد گره‌های پردازشی را از ۱ تا ۴ افزایش می‌دهیم. همان‌طور که از این اشکال مشخص است نتیجه حاصل از الگوریتم ۱ بسیار به مقدار بهینه نزدیک است.

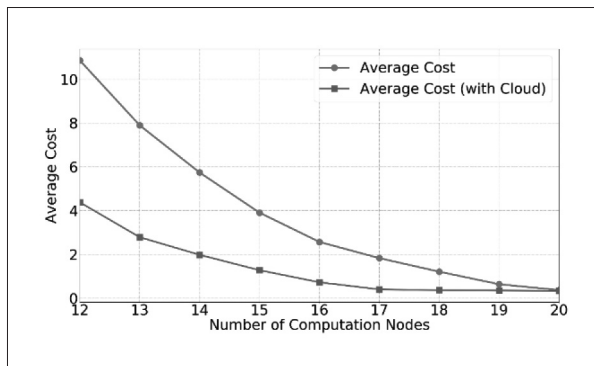
فرض می‌کنیم تعداد سرویس‌ها  $N = 15$  باشد و تعداد منابع پردازشی لبه شبکه ( $M$ ) از ۱۲ تا ۲۰ تغییر کند. دو فرآیندهای با و بدون حضور منبع پردازشی ابری را در نظر می‌گیریم و میانگین هزینه سرویس‌ها ( $-U/N$ ) را رسم می‌کنیم. انتظار داریم که با افزایش تعداد منابع پردازشی، میانگین هزینه سرویس‌ها کاهش پیدا کند. چرا که قدرت پردازشی بیش‌تری در شبکه وجود دارد و می‌تواند باعث کاهش تأخیر یا اختلاف نرخ مطلوب و نرخ بهینه شود. جهت مقایسه، موارد را با حالتی که منابع ابری نیز در شبکه وجود دارد مقایسه می‌کنیم. منابع ابری دارای قدرت پردازشی بیش‌تری هستند ولی تأخیر بالاتری دارند.

نتیجه این شبیه‌سازی در شکل ۴ آورده شده است.

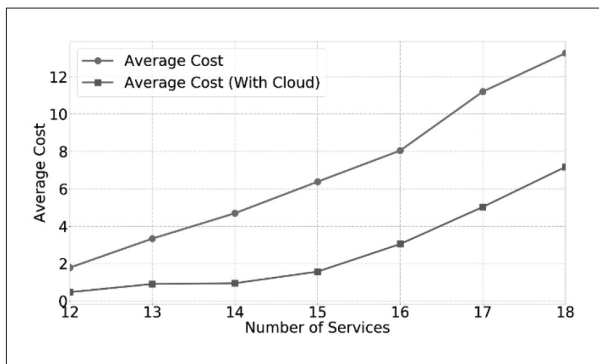




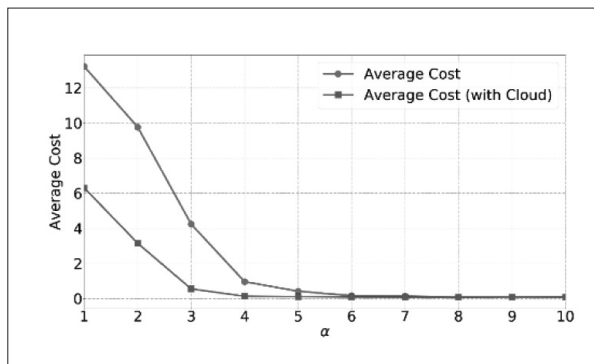
شکل ۵: تاثیر  $\omega_s$  بر تأخیر و اختلاف نرخ بهینه با نرخ مطلوب برای یک سرویس



شکل ۴: میانگین هزینه سرویس‌ها در برابر تعداد منابع پردازشی در لبه شبکه



شکل ۷: تاثیر تعداد سرویس‌ها بر میانگین هزینه سرویس‌ها



شکل ۶: تاثیر  $\alpha_s$  بر میانگین هزینه سرویس‌ها

شده رسم کرده‌ایم. در این شکل مطابق انتظار با افزایش تعداد سرویس‌ها، میانگین هزینه سرویس‌هایی که از ابتدا در سیستم قرار داشتند افزایش پیدا کرده‌است. علاوه بر این مانند بقیه شبیه‌سازی‌ها، در فرآیندهای که دارای منبع پردازشی ابری است، میانگین هزینه ۱۲ سرویس اول کمتر از سناریوی بدون منبع پردازشی ابری است. همچنین با افزایش تعداد سرویس‌ها، اختلاف این دو فرآیندها افزایش پیدا می‌کند، چراکه منبع پردازشی ابری ظرفیت پردازشی بالایی دارد و سرویس‌ها با استفاده از آن هزینه خود را کاهش می‌دهند.

#### ۵- جمع‌بندی و نتیجه‌گیری

در این مقاله مسئله تخصیص منابع پردازشی به صورت چند به چند برای سرویس‌های اینترنت اشیا بررسی شد. هر سرویس می‌تواند از چند منبع پردازشی استفاده کند و منابع پردازشی هم می‌توانند به چند سرویس اختصاص پیدا کنند. ابتدا این مسئله به صورت یک مسئله برنامه‌ریزی

نسبت اهمیت تأخیر به اختلاف نرخ بهینه و مطلوب است. در نتیجه انتظار داریم با افزایش  $\omega_s$  اختلاف نرخ، افزایش پیدا کرده و تأخیر کاهش پیدا کند که شکل هم همین را نشان می‌دهد.

در شکل ۶ تاثیر  $\alpha_s$  بر روی هزینه یک سرویس بررسی شده‌است. به دلیل تاثیر مستقیم  $\alpha_s$  در  $U_s$ ، مقدار  $-U_s/\alpha_s$  رسم شده‌است. افزایش  $\alpha_s$  باعث افزایش تاثیر  $U_s$  در تابع هدف بهینه‌سازی می‌شود. پس انتظار داریم، افزایش  $\alpha_s$  باعث کاهش  $-U_s/\alpha_s$  بشود. این کاهش معادل کاهش تأخیر و اختلاف نرخ بهینه و مطلوب سرویس  $S$  است.

در شبیه‌سازی بعدی فرض می‌کنیم که تعداد منابع پردازشی ثابت و برابر ۱۵ باشد و تعداد سرویس‌ها در بازه [12, 18] تغییر کند. به دلیل این که تعداد منابع پردازشی ثابت است، انتظار داریم افزایش تعداد سرویس‌ها باعث افزایش هزینه میانگین سرویس‌هایی که از ابتدا در سیستم وجود داشتند شود. در شکل ۷ میانگین هزینه ۱۲ سرویس اول را برای تعداد سرویس‌های شبیه‌سازی

- K. Hu, C. Chen and X. Liu, "System design of the internet of things for residential smart grid," IEEE Wireless Communications, vol. 23, no. 5, pp. 90-98, 2016.
8. R. Morabito and N. Bejar, "Enabling data processing at the network edge through lightweight virtualization technologies," in Proc. of SECON Workshops, London, United Kingdom, 2016.
9. R. Morabito, "Virtualization on internet of things edge devices with container technologies: a performance evaluation," IEEE Access, vol. 5, pp. 8835-8850, 2017.
10. R. Deng, R. Lu, C. Lai, T. H. Luan and H. Liang, "Optimal workload allocation in fog-cloud computing toward balanced delay and power consumption," IEEE Internet of Things Journal, vol. 3, no. 6, pp. 1171-1181, 2016.
11. M. Barcelo, A. Correa, J. Llorca, A. M. Tulino, J. L. Vicario and A. Morell, "IoT-cloud service optimization in next generation smart environments," IEEE JSAC, vol. 34, no. 12, pp. 4077-4090, 2016.
12. H. Zhang, Y. Xiao, S. Bu, D. Niyato, F. R. Yu and Z. Han, "Computing resource allocation in three-tier IoT fog networks: A joint optimization approach combining Stackelberg game and matching," IEEE Internet of Things Journal, vol. 4, no. 5, pp. 1204-1215, 2017.
13. L. Liu, Z. Chang, X. Guo, S. Mao and T. Ristaniemi, "Multiobjective optimization for computation offloading in fog computing," IEEE Internet of Things Journal, vol. 5, no. 1, pp. 283-295, 2017.
14. Y. Gu, Z. Chang, M. Pan, L. Song and Z. Han, "Joint radio and computational resource allocation in IoT fog computing," IEEE TVT, vol. 67, no. 8, pp. 7475-7484, 2018.
15. L. Kleinrock, "Queueing Systems," Wiley-Interscience, 1975.
16. S. Boyd and L. Vandenberghe, Convex Optimization, Cambridge university press, 2004.

غیرخطی عدد صحیح مخلوط فرمول بندی شد. با توجه به پیچیدگی حل مسئله ارائه شده، یک روش زیر بهینه برای حل آن ارائه شد. با مقایسه این روش با مقدار بهینه مشخص شد که این روش که دارای پیچیدگی محاسباتی بسیار کمتری نسبت به روش جستجوی فراگیر است، برای تعداد سرویس‌ها و گره‌های پردازشی کم به مقدار بهینه بسیار نزدیک است. نتایج شبیه‌سازی با این روش نشان می‌دهد که در این روش با افزایش تعداد گره‌های لبه شبکه، هزینه سرویس‌ها کاهش پیدا می‌کند. همچنین با استفاده از این روش، در حضور منابع پردازشی ابری، با افزایش تعداد منابع پردازشی لبه شبکه، سرویس‌های بیشتری پردازششان را به لبه شبکه انتقال می‌دهند. تا جایی که سرویس‌ها به طور کامل از منابع پردازشی لبه شبکه استفاده می‌کنند و از منابع پردازشی ابری استفاده نمی‌شود. برای کارهای آینده می‌توان روی روش‌هایی کار کرد که برای حل مسئله در ابعاد بسیار بالا با سرعت بالا به نتیجه می‌رسند یا به صورت توزیع شده جواب را به دست می‌آورند.

## مراجع

1. Cisco, "Cisco Edge-to-Enterprise IoT Analytic for Electric Utilities," 2018. [Online]. Available: <https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/big-data/solution-overview-c22-740248.pdf>.
2. W. Shi, J. Cao, Q. Zhang, Y. Li and L. Xu, "Edge computing: Vision and challenges," IEEE Internet of Things Journal, vol. 3, no. 5, pp. 637-646, 2016.
3. S. Yi, Z. Hao, Z. Qin and Q. Li, "Fog computing: Platform and applications," in Proc. of IEEE HotWeb, Washington, DC, 2015.
4. T. Pfandzelter and D. Bermbach, "tinyFaaS: A lightweight faas platform for edge environments," in 2020 IEEE International Conference on Fog Computing (ICFC), 2020.
5. K. Ha, Z. Chen, W. Hu, W. Richter, P. Pillai and M. Satyanarayanan, "Towards wearable cognitive assistance," in Proc. of international conference on Mobile systems, applications, and services (MobiSys), Bretton Woods, NH, 2014.
6. B. I. Ismail, E. M. Goortani, M. B. Ab Karim, W. M. Tat, S. Setapa, J. Y. Luke and O. H. Hoe, "Evaluation of docker as edge computing platform," in Proc. of ICOS, Melaka, Malaysia, 2015.
7. S. K. Viswanath, C. Yuen, W. Tushar, W.-T. Li, C.-K. Wen,