

# Chiến lược điều khiển và quản lý máy chủ để tăng chất lượng các dịch vụ điện toán đám mây

Trần Văn Đặng, Nguyễn Như Quỳnh

**Tóm tắt**--Trong lĩnh vực công nghiệp, theo báo cáo của Gartner [1], điện toán đám mây (Cloud Computing) là một trong mười xu hướng công nghệ của năm 2015. Mô hình điện toán này cũng đang thu hút sự quan tâm một cách mạnh mẽ cộng đồng các nhà khoa học. Cho tới nay, có rất nhiều nghiên cứu, đã và đang tập trung vào việc tiết kiệm năng lượng dành cho máy chủ trong các hệ thống đám mây. Trong quá trình vận hành của các trung tâm dữ liệu đám mây, các máy chủ không tải (idle) tiêu tốn một lượng lớn năng lượng không cần thiết. Để khắc phục vấn đề này, nhiều giải pháp được đề xuất như tắt các máy chủ khi không thực hiện công việc nào. Tuy nhiên việc tắt bật các máy chủ ảnh hưởng không nhỏ đến chất lượng dịch vụ điện toán đám mây do việc khởi động các máy thường tiêu tốn thời gian nhất định làm giảm khả năng phục vụ nhanh chóng của các dịch vụ chạy trên đó. Trong bài báo này, chúng tôi đưa ra chiến lược quản lý máy chủ và cung cấp tài nguyên một cách hiệu quả, làm giảm thời gian chờ đợi phục vụ của các dịch vụ, theo cách đó tăng chất lượng dịch vụ điện toán đám mây. Chúng tôi sử dụng công cụ CloudSIM[8] để cài đặt và đánh giá chiến lược đã đề xuất. Kết quả thu được đã chứng minh hiệu quả của mô hình chúng tôi đưa ra dưới các điều kiện khác nhau so với các mô hình trước đó.

**Từ khóa**--điện toán đám mây, điện toán xanh, lý thuyết hàng đợi, chuỗi Markov.

## 1. GIỚI THIỆU

Điện toán đám mây là một trong các xu hướng phát triển mạnh mẽ nhất những năm gần đây. Mô hình này đang thu hút sự quan tâm nghiên cứu của nhiều nhà khoa học. Động lực của sự hấp dẫn này chính là nhờ những ưu điểm mà đám mây mang lại cho người sử dụng như: sự mềm dẻo (elasticity), tính khả mở (scalability), tính sẵn sàng cao (availability) của hệ thống, cộng với mô hình trả tiền theo sự tiêu dùng tài nguyên thực tế (pay-as-you-go) giúp tăng hiệu quả và tiết kiệm chi phí tối đa cho người dùng. Đi đôi với sự phát triển và mở rộng của các dịch vụ điện toán đám mây, các nhà cung cấp thường dành mối quan tâm lớn để đảm bảo chất lượng dịch vụ (Quality of Service – QoS) của mình. Chỉ có như vậy, họ (các nhà cung cấp) mới có thể tìm kiếm được lợi nhuận trong thị phần rất nhiều các dịch vụ đám mây ngày nay.

Mức tiêu thụ năng lượng của các trung tâm dữ liệu là vấn đề quan trọng đối với các nhà cung cấp dịch vụ điện toán đám mây.

Công trình này được thực hiện dưới sự hướng dẫn của TS Nguyễn Bình Minh.

Trần Văn Đặng, sinh viên lớp KSTN - CNTT, khóa 57, Viện Công nghệ thông tin và Truyền thông, trường Đại học Bách Khoa Hà Nội (điện thoại: 01626489218, e-mail: dangtran.hust@gmail.com).

Nguyễn Như Quỳnh, sinh viên lớp KSCLC-HTTT&TT, khóa 56, Viện Công nghệ thông tin và Truyền thông, trường Đại học Bách Khoa Hà Nội (điện thoại: 01657645549, e-mail: nguyen.nhu.quynh.1993@gmail.com).

© Viện Công nghệ thông tin và Truyền thông, trường Đại học Bách Khoa Hà Nội.

Cho tới nay, có rất nhiều nghiên cứu, đã và đang tập trung vào việc tiết kiệm năng lượng dành cho máy chủ trong các hệ thống đám mây [2], [3]. Saiqui Long và đồng nghiệp [4] đưa ra mô hình ba bước tiết kiệm năng lượng cho các kho lưu trữ đám mây dựa trên phân phối Poisson, độ dài hàng đợi và chuỗi Markov. Một vài nghiên cứu định nghĩa và tập trung hẳn sang xu hướng điện toán xanh (green computing) [5], [6] chủ yếu đưa ra việc quản lý hiệu quả quá trình bật và tắt các máy chủ vật lý sao cho tối ưu nhất năng lượng tiêu thụ máy chủ. Trong bài báo của mình, Phung-Duc Tuan [7] cũng sử dụng phân phối Poisson và chuỗi Markov để chứng minh được khả năng phục vụ tối ưu các công việc (jobs) và sự tiêu thụ năng lượng của hệ thống máy chủ với hai trạng thái của máy được định nghĩa là OFF (tắt máy) và ON (đang phục vụ công việc). Ở đây, tác giả giả sử rằng, các công việc sau khi đã hoàn thành thì máy đang ở trạng thái ON sẽ tự động chuyển sang OFF. Khi có công việc tới, tương tự như vậy, các máy chủ sẽ từ trạng thái OFF chuyển sang ON. Tuy nhiên, trong thực tế, hầu hết các nhà cung cấp dịch vụ điện toán đám mây đều luôn luôn phải giữ một số lượng máy chủ ở một trạng thái “trung gian”, khi đó máy chủ được bật nhưng chưa nhận bất kỳ một công việc nào. Với mô hình ba trạng thái (tắt, trung gian và bật), khả năng đáp ứng nhu cầu (trong thực tế luôn thay đổi nhanh chóng) sử dụng tài nguyên của các dịch vụ/ứng dụng trong môi trường đám mây mới được đảm bảo nhờ tốc độ đáp ứng công việc nhanh chóng. Theo cách đó, việc giữ số lượng máy ở trạng thái trung gian đồng thời sẽ giúp tăng chất lượng dịch vụ (QoS) mà nhà cung cấp dịch vụ điện toán đám mây phân phối cho người sử dụng.

Trong bài báo này, ngoài hai trạng thái ON và OFF chúng tôi định nghĩa một trạng thái trung gian (MIDDLE) cho các máy chủ vật lý trong trung tâm dữ liệu đám mây. Với trạng thái MIDDLE này, một số lượng máy chủ nhất định được bật và chờ đợi các công việc tới để phục vụ. Cũng sử dụng phân phối Poisson và chuỗi Markov, chúng tôi đưa ra các thuật toán xác định khi nào thì cần bật thêm máy từ OFF sang MIDDLE, nhằm đảm bảo khả năng phục vụ công việc của hệ thống, từ đó cho phép tăng chất lượng dịch vụ đám mây. Mô hình ba trạng thái và các thuật toán điều khiển bật, tắt máy chủ của chúng tôi được cài đặt và thực nghiệm thành công dựa trên phần mềm mô phỏng môi trường đám mây CloudSIM [8]. Các kết quả thu được thông qua thí nghiệm mô phỏng đã chứng minh hiệu quả đạt được của mô hình đề xuất mới này.

## 2. CÁC CÔNG TRÌNH LIÊN QUAN

Hiện nay các chiến lược cung cấp tài nguyên phù hợp trong môi trường điện toán đám mây đang được quan tâm và nghiên cứu khi mà các trung tâm dữ liệu tiêu thụ một lượng lớn năng

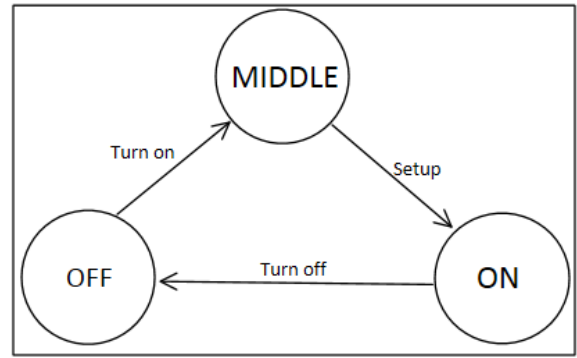
lượng hàng năm. Mặc dù với mô hình hai trạng thái ON, OFF của các máy chủ được đưa ra trong các bài báo [7, 9, 10] các tác giả chứng minh được công thức tính chính xác cho thời gian chờ đợi của công việc vào hệ thống và tính toán năng lượng tiêu thụ, tuy nhiên với mô hình hai trạng thái này, thời gian chờ đợi phục vụ công việc luôn lớn do các máy ON được chuyển về OFF nếu không thực hiện công việc nào hoặc máy từ OFF chuyển sang ON tốn khá nhiều thời gian khởi động, cài đặt và cấu hình ứng dụng. Trong thực tế, để đảm bảo tốc độ phục vụ nhanh hầu hết các nhà cung cấp dịch vụ điện toán đám mây đều luôn luôn phải giữ một số lượng máy chủ để sẵn sàng thực hiện khi có công việc mới. Số lượng máy chủ giữ ở trạng thái trung gian phải phù hợp với cường độ công việc, cũng như khả năng phục vụ của hệ thống, đảm bảo giảm được thời gian chờ đợi mà không lãng phí điện năng cho việc duy trì lượng máy chủ đó. Một vài công trình đưa ra các thuật toán quyết định việc bật tắt các máy chủ ở trạng thái không tải [4] theo cường độ công việc, và tập trung vào việc tiết kiệm năng lượng dành cho máy chủ trong các hệ thống đám mây [2], [3], hay quản lý hiệu quả quá trình bật và tắt các máy chủ vật lý [5], [6]. Trong công trình của chúng tôi, mô hình ba trạng thái ON, OFF, MIDDLE được đề xuất, số lượng các máy chủ cần bật được tính toán một cách phù hợp với nhu cầu công việc và khả năng của hệ thống. Đóng góp trong nghiên cứu của chúng tôi cụ thể là:

- Đưa ra mô hình ba trạng thái cho máy chủ trong các trung tâm dữ liệu của đám mây.
- Chứng minh mô hình dựa trên lý thuyết kết hợp lý thuyết hàng đợi và chuỗi (xích) Markov.
- Mô phỏng mô hình đề xuất bằng công cụ CloudSim, chứng minh sự hiệu quả trong việc nâng cao chất lượng dịch vụ các dịch vụ chạy trên máy chủ đám mây.

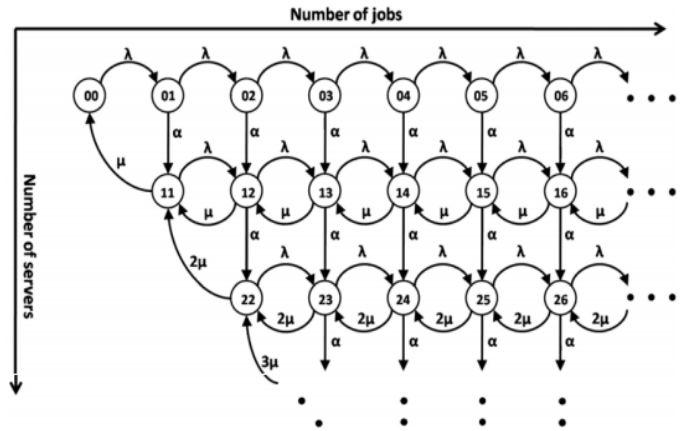
### 3. MÔ HÌNH LÝ THUYẾT

#### 3.1. Mô hình ON/OFF/MIDDLE/ $\infty$ /STAG

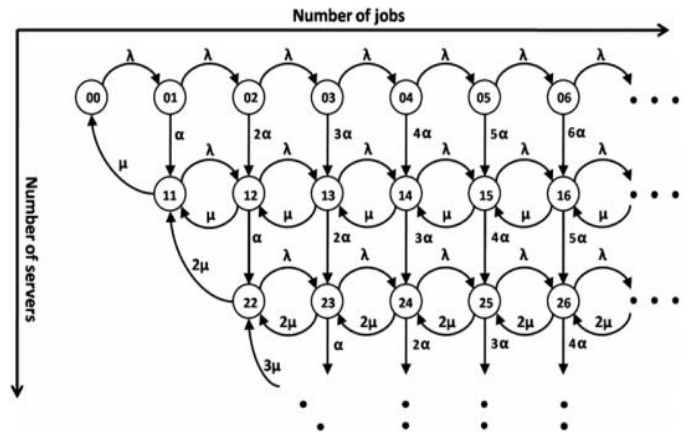
Trong các trung tâm dữ liệu, thường các máy chủ vật lý không tải (không phục vụ công việc – job nào) sẽ được tắt đi để tiết kiệm điện năng. Khi một máy chủ được bật lại, để phục vụ được các công việc máy chủ sẽ tiêu tốn chi phí thiết lập lại hệ thống bao gồm cả thời gian cài đặt và năng lượng [9]. Thực tế, các trung tâm dữ liệu như: Google, Microsoft, Yahoo và Amazon chứa hàng chục ngàn máy chủ, do vậy có thể coi số lượng máy chủ là không giới hạn. Trong bài [9] Gandhi đưa ra mô hình ON/OFF/ $\infty$ /STAG cho các hệ thống đó. Như đã nhắc tới ở phần 1, chúng tôi định nghĩa thêm một trạng thái trung gian MIDDLE cho các máy chủ vật lý trong trung tâm dữ liệu đám mây. Khi một máy chủ MIDDLE được chuyển sang trạng thái ON nó phải được đặt vào chế độ Setup (hình 3.1). Trong quá trình đó, máy chủ không thể phục vụ bất kỳ công việc nào. Thời gian để đưa một máy chủ từ chế độ setup sang trạng thái ON gọi là thời gian SETUP (SETUP time). Chúng tôi giả sử thời gian đó là biến I, có phân phối mũ với  $E[I] = 1/\alpha$ . Chúng tôi cũng giả sử các công việc (job) đến hệ thống theo quá trình Poisson với tốc độ  $\lambda$  và thời gian để hoàn thành một công việc là X có phân phối mũ với  $E[X] = 1/\mu$ . Khi một máy chủ không được sử dụng nó được tắt về trạng thái OFF ngay lập tức. Khi một công việc mới vào hệ thống, công việc này được đưa vào hàng đợi, chờ một máy chủ



Hình 3.1: Sơ đồ chuyển đổi trạng thái của máy chủ vật lý



Hình 3.2: Chuỗi markov với ON/OFF/MIDDLE/ $\infty$ /STAG



Hình 3.3: Chuỗi markov với ON/OFF/MIDDLE/ $\infty$

đang trong quá trình SETUP chuyển sang ON để thực hiện công việc này, trường hợp không có, một máy ở trạng thái MIDDLE (nếu có) sẽ được chuyển sang chế độ SETUP. Khi một máy chủ  $j$  hoàn thành công việc, công việc đầu tiên trong hàng đợi được chuyển sang máy chủ  $j$  mặc dù công việc này đang đợi một máy chủ  $i$  trong trạng thái SETUP, sau đó máy chủ  $i$  sẽ được tắt. Trong mô hình này chúng tôi chỉ cho phép nhiều nhất một máy chủ trong chế độ SETUP tại bất kỳ thời điểm nào, chính sách này được nhắc đến là stagger setup trong bài [9]. Các máy OFF được bật sang MIDDLE theo một thuật toán kiểm soát mà chúng tôi đề xuất để đảm bảo luôn có máy chủ ở MIDDLE khi hệ thống cần. Số lượng máy ở MIDDLE cũng không được phép quá lớn

tránh lãng phí điện năng. Quá trình máy chủ chuyển đổi từ OFF sang MIDDLE mất một thời gian như nhau đối với các máy chủ giống nhau ( $t_{\text{OFF} \rightarrow \text{MIDDLE}}$ ). Giả sử luôn có máy chủ ở trạng thái MIDDLE, ta có thể biểu diễn sự chuyển đổi trạng thái của hệ thống bằng xích markov (markov chain) ON/OFF/ $\infty$ /STAG [9] (hình 3.2) với trạng thái  $(i, j)$ , trong đó  $i$  là số lượng máy chủ ON,  $j$  là số lượng công việc trong hệ thống. Khi  $j > i$  ta có 1 máy chủ trong SETUP và khi  $i = j$ , không có máy chủ nào trong SETUP.

### 3.2. Mô hình ON/OFF/MIDDLE/ $\infty$

Mô hình này tương tự với mô hình ON/OFF/MIDDLE/ $\infty$ /STAG ngoại trừ việc cho phép nhiều máy chủ đồng thời có thể trong chế độ SETUP, sự chuyển đổi trạng thái của hệ thống được biểu diễn trong hình 3.3.

## 4. MÔ HÌNH VÀ THUẬT TOÁN CUNG CẤP TÀI NGUYÊN ĐỀ XUẤT

Như đã đề xuất trong phần 3, cần phải điều khiển việc bật các máy chủ OFF sang MIDDLE sao cho trong hệ thống luôn có máy chủ ở trạng thái MIDDLE mà số máy chủ này không quá lớn. Để đảm bảo cân bằng giữa hai điều kiện trên, chúng tôi đưa ra con số trung bình cho số lượng máy cần bật từ OFF sang MIDDLE sau một khoảng thời gian cho trước.

**Bổ đề 1:** Xét  $X_1, X_2, \dots, X_n$  là các biến ngẫu nhiên có phân phối mũ với kỳ vọng  $\mu_1^{-1}, \mu_2^{-1}, \dots, \mu_n^{-1}$  là và chúng độc lập với nhau. Xét  $X = \min \{X_1, X_2, \dots, X_n\}$ . Khi đó  $X$  cũng có phân phối mũ với kỳ vọng  $(\mu_1 + \mu_2 + \dots + \mu_n)^{-1}$

Chứng minh:

$$\begin{aligned} P(X > x) &= P(\min \{X_1, X_2, \dots, X_n\} > x) \\ &= P(X_1 > x) \cdot P(X_2 > x) \dots P(X_n > x) \\ &= (e^{-\mu_1 x}) \cdot (e^{-\mu_2 x}) \dots (e^{-\mu_n x}) \\ &= e^{-(\mu_1 + \mu_2 + \dots + \mu_n)x} \end{aligned}$$

Suy ra  $P(X < x) = 1 - e^{-(\mu_1 + \mu_2 + \dots + \mu_n)x}$  (phân phối mũ)

**Bổ đề 2:** Giả sử  $X_1, X_2$  là hai biến ngẫu nhiên độc lập có phân phối mũ với kỳ vọng tương ứng là  $\mu_1^{-1}, \mu_2^{-1}$ . Khi đó:

$$P(X_1 < X_2) = \frac{\mu_1}{\mu_1 + \mu_2}$$

Chứng minh: Theo công thức xác suất đầy đủ ta có:

$$\begin{aligned} P(X_1 < X_2) &= \int_0^{+\infty} P(X_1 < X_2 | X_2 = x) \mu_2 e^{-\mu_2 x} dx \\ &= \int_0^{+\infty} P(X_1 < x) \mu_2 e^{-\mu_2 x} dx = \frac{\mu_1}{\mu_1 + \mu_2} \end{aligned}$$

**Bổ đề 3:** Xét  $X_1, X_2, \dots, X_n$  là các biến ngẫu nhiên có phân phối mũ với trung bình là  $\mu_1^{-1}, \mu_2^{-1}, \dots, \mu_n^{-1}$  và chúng độc lập với nhau. Khi đó

$$P\{X_i = \min(X_1, X_2, \dots, X_n)\} = \frac{\mu_i}{\mu_1 + \mu_2 + \dots + \mu_n}$$

Đặt  $y = \min(X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ . Theo bổ đề 1,  $y$  có phân phối mũ với kỳ vọng là:

$$(\mu_1 + \mu_2 + \dots + \mu_{i-1} + \mu_{i+1} + \dots + \mu_n)^{-1}. \text{ Suy ra: } P\{X_i = \min(X_1, X_2, \dots, X_n)\}$$

$$= P(X_i < y) = \frac{\mu_i}{\mu_1 + \mu_2 + \dots + \mu_n} \quad (\text{theo bổ đề 2})$$

**Bổ đề 4:** giả sử  $X_1, X_2$  là hai biến ngẫu nhiên độc lập có phân phối mũ với trung bình tương ứng là  $\mu_1^{-1}, \mu_2^{-1}$ . Khi đó:

$$P\{\min(X_1, X_2) > t | X_1 < X_2\} = P\{\min(X_1, X_2) > t\}$$

Hay nói cách khác xác suất của  $\min(X_1, X_2)$  và xác suất để  $\min(X_1, X_2)$  là  $X_1$  hay  $X_2$  là độc lập

Chứng minh:

$$P\{\min(X_1, X_2) > t\} = e^{-(\mu_1 + \mu_2)t}$$

$$P\{\min(X_1, X_2) > t | X_1 < X_2\} = P(X_1 > t | X_1 < X_2)$$

$$= \frac{P(X_1 > t, X_1 < X_2)}{P(X_1 < X_2)} = \frac{P(t < X_1 < X_2)}{P(X_1 < X_2)}$$

Theo công thức xác suất đầy đủ:

$$\begin{aligned} P(t < X_1 < X_2) &= \int_t^{+\infty} P(t < X_1 < X_2 | X_2 = x) \mu_2 e^{-\mu_2 x} dx \\ &= \int_t^{+\infty} P(t < X_1 < x) \mu_2 e^{-\mu_2 x} dx = \int_t^{+\infty} (e^{-\mu_1 t} - e^{-\mu_1 x}) \mu_2 e^{-\mu_2 x} dx \\ &= \frac{\mu_1}{\mu_1 + \mu_2} e^{-(\mu_1 + \mu_2)t} \end{aligned}$$

$$P(X_1 < X_2) = \frac{\mu_1}{\mu_1 + \mu_2}$$

Suy ra  $P\{\min(X_1, X_2) > t | X_1 < X_2\} = P\{\min(X_1, X_2) > t\}$

### 4.1 Thuật toán với ON/OFF/MIDDLE/ $\infty$ /STAG

Khi luôn có máy chủ ở trạng thái MIDDLE, sự chuyển đổi trạng thái của hệ thống được thể hiện bằng xích markov ON/OFF/ $\infty$ /STAG (hình 3.2). Trong bài [9], Gandhi và Harchol-Balter chứng minh xác suất giới hạn cho xích markov ON/OFF/ $\infty$ /STAG được tính bởi công thức:

$$\pi_{i,j} = \frac{\pi_{0,0} \rho^i}{i!} \left( \frac{\lambda}{\lambda + \alpha} \right)^{j-i} \quad \text{và} \quad \pi_{0,0} = \sum_{i=0}^k \sum_{j \geq i} \frac{\rho^i}{i!} \left( \frac{\lambda}{\lambda + \alpha} \right)^{j-i}$$

**Bổ đề 5:** Với mô hình ON/OFF/MIDDLE/ $\infty$ /STAG. Xét tại thời điểm  $t$ , hệ thống đang ở trạng thái  $(\theta, i)$ ,  $i > 0$ . Tại thời điểm  $t + h$  xác suất để hệ thống chuyển sang trạng thái mới  $(\theta', j)$  ( $\theta' \neq \theta$  hoặc  $i \neq j$ ) là:

$$P\{(\theta, i) \rightarrow (\theta', j)\} = \begin{cases} \alpha h + o(h), & \theta' = \theta + 1 \text{ và } i = j \\ \lambda h + o(h), & \theta' = \theta \text{ và } j = i + 1 \\ \theta \mu h + o(h), & \theta' = \theta \text{ và } j = i - 1 \\ o(h), & \text{trường hợp khác} \end{cases}$$

Chứng minh: Khi  $i > 0$ , trong hàng đợi có công việc (job). Như vậy đến thời điểm  $t + h$ , có thể có các sự kiện xảy ra: có job đến, có job hoàn thành, có thêm máy chủ chuyển sang ON từ chế độ SETUP

Theo bổ đề 1, thời gian từ  $t$  đến sự kiện tiếp theo (có 1 job đến hoặc có 1 job hoàn thành, hoặc có một máy chuyển sang ON) cũng có phân phối mũ với kỳ vọng là  $(\lambda + \alpha + \theta\mu)^{-1}$  (vì thời gian đến sự kiện tiếp theo là giá trị nhỏ nhất trong: thời gian job tiếp theo đến (interarrival time),  $\theta$  service time (thời gian hoàn thành công việc) của  $\theta$  server ON, SETUP time; các giá trị này đều là các biến ngẫu nhiên có phân phối mũ).

Do đó xác suất để sự kiện xảy ra trong khoảng  $(t, t+h)$  là:

$$1 - e^{-(\lambda + \alpha + \theta\mu)h} = (\lambda + \alpha + \theta\mu)h + o(h)$$

Theo bổ đề 3, xác suất để sự kiện xảy ra là job đến:  $\alpha/(\lambda + \alpha + \theta\mu)$ . Mặt khác thời gian để xảy ra sự kiện tiếp theo và loại của sự kiện đó là 2 biến độc lập (bổ đề 4) nên xác suất để hệ thống chuyển trạng thái từ  $(\theta, i) \rightarrow (\theta, i+1)$  là:

$$P\{(\theta, i) \rightarrow (\theta, i+1)\}$$

$$= [(\lambda + \alpha + \theta\mu)h + o(h)] \frac{\lambda}{\lambda + \alpha + \theta\mu} + o(h) = \lambda h + o(h)$$

Xác suất sự kiện xảy ra là job hoàn thành là  $\theta\mu/(\lambda + \alpha + \theta\mu)$  nên:

$$P\{(\theta, i) \rightarrow (\theta, i-1)\}$$

$$= [(\lambda + \alpha + \theta\mu)h + o(h)] \frac{\theta\mu}{\lambda + \alpha + \theta\mu} + o(h) = \theta\mu h + o(h)$$

Tương tự:

$$P\{(\theta, i) \rightarrow (\theta+1, i)\}$$

$$= [(\lambda + \alpha + \theta\mu)h + o(h)] \frac{\alpha}{\lambda + \alpha + \theta\mu} + o(h) = \alpha h + o(h)$$

**Định lý 1:** Với mô hình ON/OFF/MIDDLE/ $\infty$ /STAG, số lượng máy chủ ở trạng thái MIDDLE và SETUP trung bình giảm đi trong thời gian  $t$  là:

$$E[K] = \frac{\lambda \alpha t}{\lambda + \alpha} \quad (1)$$

Chứng minh: gọi  $K_h$  là số máy trong trạng thái MIDDLE hoặc SETUP giảm đi trong thời gian  $h$  nhỏ. Lượng máy chủ ở MIDDLE và SETUP chỉ giảm đi khi có một máy chủ mới được bật sang ON. Do vậy  $K_h > 0$  khi hệ thống chuyển trạng thái từ  $(\theta, i) \rightarrow (\theta', j)$  với  $\theta' > \theta$ . Từ xích markov và bổ đề 5,  $P(K_h > 1) = o(h)$ .  $K_h = 1$  khi hệ chuyển trạng thái từ  $(\theta, i) \rightarrow (\theta+1, i)$ , với  $i > 0$ . Do đó, theo công thức xác suất đầy đủ:

$$P(K_h = 1) = \sum_{i>0} \pi_{\theta,i} \cdot P\{(\theta, i) \rightarrow (\theta+1, i)\}$$

$$= \left(1 - \sum_{i=0}^{\infty} \pi_{i,i}\right) \cdot (\alpha h + o(h))$$

$$= (\alpha h + o(h)) \left(1 - \sum_{i=0}^{\infty} \frac{\pi_{0,0} \cdot \rho^i}{i!}\right)$$

$$= (\alpha h + o(h)) \left(1 - \pi_{0,0} \cdot \sum_{i=0}^{\infty} \frac{\rho^i}{i!}\right)$$

$$= (\alpha h + o(h)) \left(1 - \frac{\sum_{i=0}^{\infty} \frac{\rho^i}{i!}}{\sum_{i=0}^k \sum_{j \geq i} \frac{\rho^i}{i!} \left(\frac{\lambda}{\lambda + \alpha}\right)^{j-i}}\right)$$

$$= (\alpha h + o(h)) \left(1 - \frac{\sum_{i=0}^{\infty} \frac{\rho^i}{i!}}{\sum_{i=0}^k \sum_{l \geq 0} \frac{\rho^i}{i!} \left(\frac{\lambda}{\lambda + \alpha}\right)^l}\right)$$

$$= (\alpha h + o(h)) \left(1 - \frac{\sum_{i=0}^{\infty} \frac{\rho^i}{i!}}{\sum_{i=0}^k \frac{\rho^i}{i!} \sum_{l \geq 0} \left(\frac{\lambda}{\lambda + \alpha}\right)^l}\right)$$

$$= (\alpha h + o(h)) \left(1 - \frac{1}{\sum_{l \geq 0} \left(\frac{\lambda}{\lambda + \alpha}\right)^l}\right)$$

$$= (\alpha h + o(h)) \left(1 - \frac{1}{1 / \left(1 - \frac{\lambda}{\lambda + \alpha}\right)}\right) = \frac{\lambda}{\lambda + \alpha} (\alpha h + o(h))$$

$$\text{Suy ra } E[K_h] = \frac{\lambda}{\lambda + \alpha} (\alpha h + o(h))$$

Chia  $t$  thành các khoảng thời gian  $h$  đủ nhỏ gọi  $K_i$  là số lượng máy MIDDLE và SETUP giảm trong khoảng thời gian thứ  $i$ . Ta có:

$$E[K_i] = \frac{\lambda}{\lambda + \alpha} (\alpha h + o(h)) \approx \frac{\lambda}{\lambda + \alpha} \alpha h. \quad K = \sum K_i. \quad \text{Do đó:}$$

$$E[K] = E[\sum K_i] = \sum E[K_i] = \frac{\lambda}{\lambda + \alpha} \alpha h \cdot \frac{t}{h} = \frac{\lambda \alpha t}{\lambda + \alpha}$$

Số lượng máy chủ trung bình giảm đi sau thời gian  $t$  được tính theo công thức (1), do vậy để bù lại số lượng giảm đi đó cần bật máy từ OFF sang MIDDLE. Khoảng thời gian trung bình để bật thêm một máy:

$$\tau = \frac{t}{E[K]} = \frac{\lambda + \alpha}{\lambda \alpha} \quad (2)$$

Sau đây là thuật toán điều khiển bật máy từ OFF sang MIDDLE:

**Thuật toán 1.** Thuật toán bật máy với mô hình ON/OFF/MIDDLE/ $\infty$ /STAG

- 1: **while** true **do**
- 2:   bật 1 máy chủ từ OFF sang MIDDLE
- 3:   tính thời gian  $\tau$  theo công thức (2)
- 4:   chờ khoảng thời gian  $\tau$
- 5: **end while**

#### 4.2. Thuật toán với ON/OFF/MIDDLE/ $\infty$

Trong mô hình này, nhiều máy chủ có thể đồng thời ở trạng thái SETUP. Sự chuyển đổi trạng thái của hệ thống khi luôn có máy chủ ở trạng thái MIDDLE được thể hiện bằng xích markov ON/OFF/ $\infty$  (hình 3.3). Xác suất giới hạn cho xích markov ON/OFF/ $\infty$  được tính bởi công thức [9]:

$$\pi_{i,j} = \frac{\pi_{0,0} \cdot \rho^i}{i!} \prod_{l=1}^{j-i} \frac{\lambda}{\lambda + l\alpha}, \quad i \geq 0, j \geq i$$

$$\pi_{0,0} = e^{-\rho} \left( \sum_{j=0}^{\infty} \prod_{l=1}^j \frac{\lambda}{\lambda + l\alpha} \right)^{-1}$$

**Định lý 2:** với mô hình ON/OFF/MIDDLE/ $\infty$ , số lượng máy chủ ở trạng thái MIDDLE và SETUP trung bình giảm đi trong thời gian  $t$  là:

$$E[K] = \alpha t \cdot \frac{\sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}{\sum_{k=0}^{\infty} \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}} \quad (3)$$

Chúng minh: Hoàn toàn tương tự với mô hình ON/OFF/MIDDLE/ $\infty$ /STAG:

$$P(K_h = 1) = \sum_{j>i} \pi_{i,j} \cdot P\{(i, j) \rightarrow (i+1, j)\}$$

Áp dụng bổ đề 5 đối với mô hình ON/OFF/MIDDLE/ $\infty$

$$P\{(i, j) \rightarrow (i+1, j)\} = (j-i) \cdot \alpha h + o(h)$$

Suy ra:

$$P(K_h = 1) = \sum_{j>i} \frac{\pi_{0,0} \cdot \rho^i}{i!} \prod_{l=1}^{j-i} \frac{\lambda}{\lambda + l\alpha} \cdot ((j-i) \cdot \alpha h + o(h))$$

$$\approx \pi_{0,0} \cdot \sum_{j>i} \frac{\rho^i}{i!} \prod_{l=1}^{j-i} \frac{\lambda}{\lambda + l\alpha} \cdot ((j-i) \cdot \alpha h)$$

$$= \pi_{0,0} \cdot \sum_{i=0}^{\infty} \frac{\rho^i}{i!} \sum_{j=i}^{\infty} \prod_{l=1}^{j-i} \frac{\lambda}{\lambda + l\alpha} \cdot ((j-i) \cdot \alpha h)$$

$$= \pi_{0,0} \cdot \alpha h \cdot \sum_{i=0}^{\infty} \frac{\rho^i}{i!} \sum_{k=0}^{\infty} \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}$$

$$= \pi_{0,0} \cdot \alpha h \cdot \left( \sum_{i=0}^{\infty} \frac{\rho^i}{i!} \right) \sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}$$

$$= \alpha h \cdot \frac{\left( \sum_{i=0}^{\infty} \frac{\rho^i}{i!} \right) \sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}{e^{\rho} \left( \sum_{j=0}^{\infty} \prod_{l=1}^j \frac{\lambda}{\lambda + l\alpha} \right)}$$

$$= \alpha h \cdot \frac{\sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}{\sum_{k=0}^{\infty} \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}$$

Chia  $t$  thành các khoảng thời gian  $h$  đủ nhỏ gọi  $K_i$  là số lượng máy MIDDLE và SETUP giảm trong khoảng thời gian thứ  $i$ . ta có:

$$E[K] = \sum E[K_i] = \frac{t}{h} \alpha h \cdot \frac{\sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}{\sum_{k=0}^{\infty} \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}$$

$$\begin{aligned} & \sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha} \\ &= \alpha t \cdot \frac{\sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}{\sum_{k=0}^{\infty} \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}} \end{aligned}$$

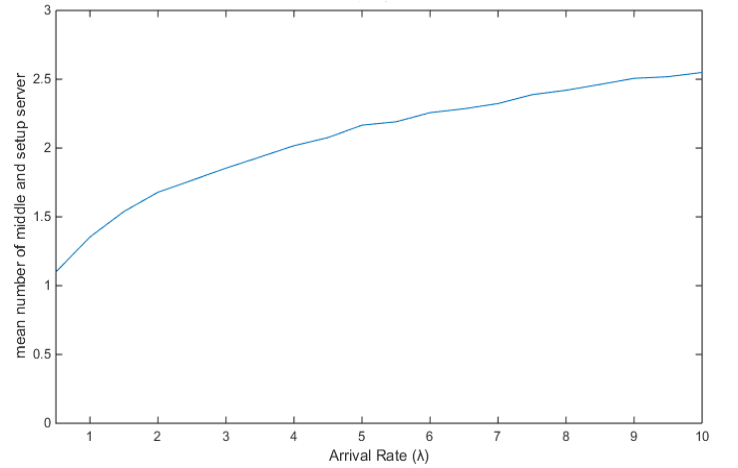
Khoảng thời gian để bật máy OFF sang MIDDLE đối với mô hình ON/OFF/MIDDLE/ $\infty$

$$\tau = \frac{t}{E[K]} = \frac{1}{\alpha} \cdot \frac{\sum_{k=0}^{\infty} \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}}{\sum_{k=0}^{\infty} k \cdot \prod_{l=1}^k \frac{\lambda}{\lambda + l\alpha}} \quad (4)$$

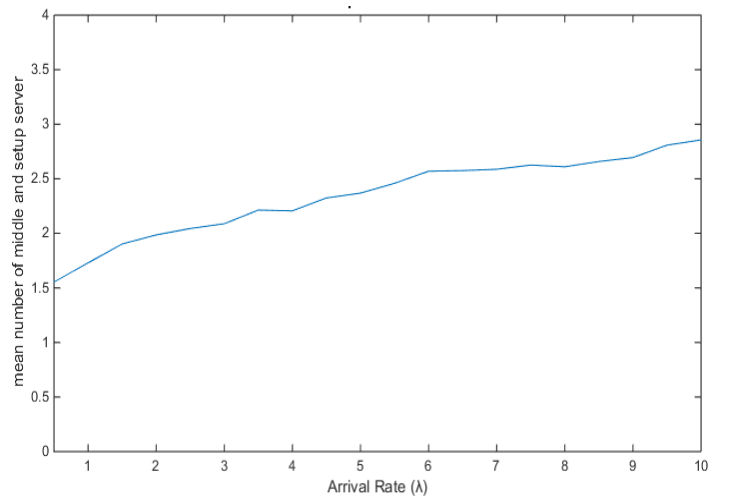
Ta có thuật toán bật máy thứ hai áp dụng với mô hình ON/OFF/MIDDLE/ $\infty$

#### Thuật toán 2

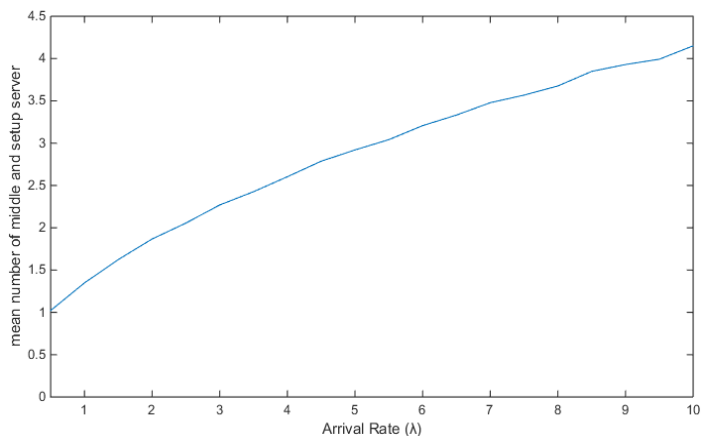
- 1: **while** true **do**
- 2:   bật 1 máy chủ từ OFF sang MIDDLE
- 3:   tính thời gian  $\tau$  theo công thức (4)
- 4:   chờ khoảng thời gian  $\tau$
- 5: **end while**



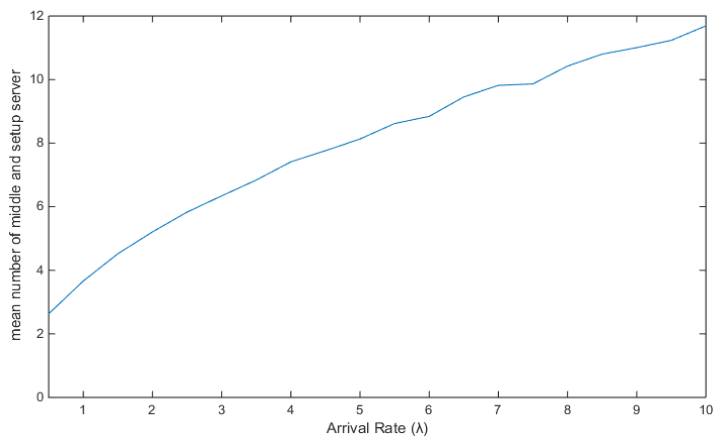
Hình 5.1a: Số máy chủ MIDDLE và SETUP trung bình mô hình ON/OFF/MIDDLE/ $\infty$ /STAG,  $\alpha = 0.1$



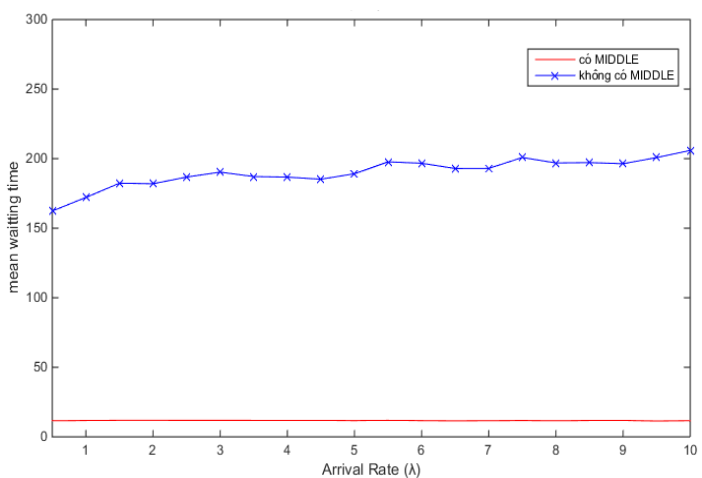
Hình 5.1b: Số máy chủ MIDDLE và SETUP trung bình mô hình ON/OFF/MIDDLE/ $\infty$ /STAG,  $\alpha = 0.01$



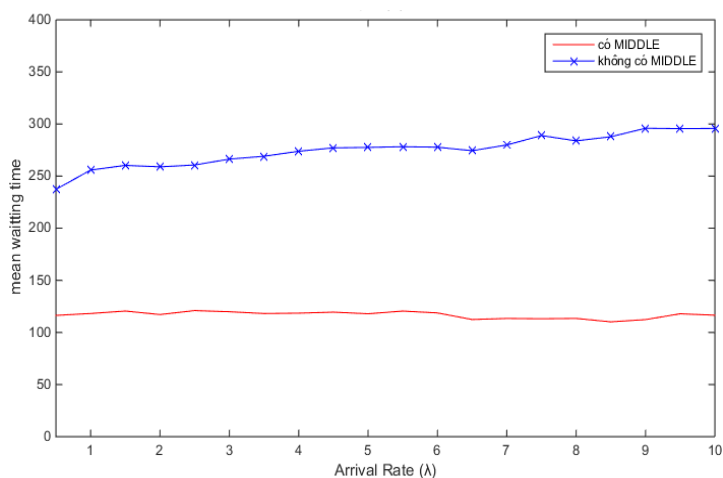
Hình 5.2a: Số máy chủ MIDDLE và SETUP trung bình mô hình ON/OFF/MIDDLE/ $\infty$ ,  $\alpha = 0.1$



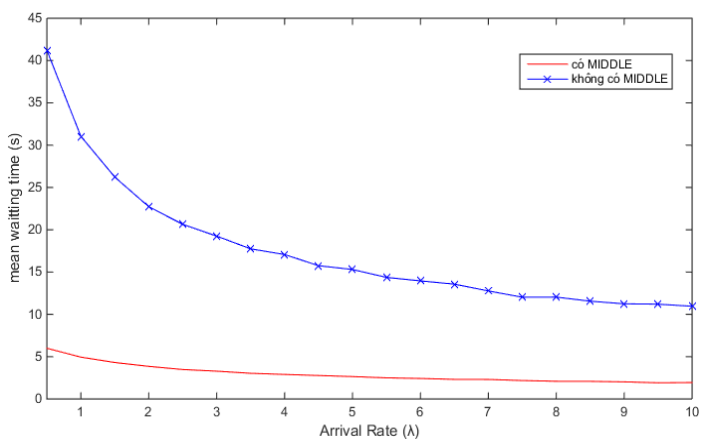
Hình 5.2b: Số máy chủ MIDDLE và SETUP trung bình mô hình ON/OFF/MIDDLE/ $\infty$ ,  $\alpha = 0.01$



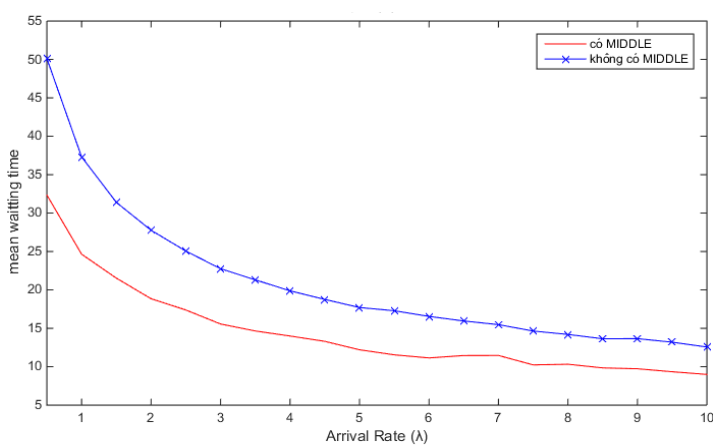
Hình 5.3a: Thời gian chờ đợi trung bình mô hình ON/OFF/MIDDLE/ $\infty$ /STAG,  $\alpha = 0.1$



Hình 5.3b: Thời gian chờ đợi trung bình mô hình ON/OFF/MIDDLE/ $\infty$ /STAG,  $\alpha = 0.01$



Hình 5.4a: Thời gian chờ đợi trung bình mô hình ON/OFF/MIDDLE/ $\infty$ ,  $\alpha = 0.1$



Hình 5.4b: Thời gian chờ đợi trung bình mô hình ON/OFF/MIDDLE/ $\infty$ ,  $\alpha = 0.01$

## 5. THỬ NGHIỆM VÀ ĐÁNH GIÁ

Trong phần này, chúng tôi đánh giá thuật toán của mình qua thời gian chờ đợi và lượng máy MIDDLE trung bình có trong hệ thống. Chúng tôi sử dụng công cụ CloudSim [8] cho việc mô phỏng các trung tâm dữ liệu, các máy chủ vật lý cũng như các công việc đến hệ thống; sử dụng thuật toán sinh ngẫu nhiên để sinh ra các giá trị ngẫu nhiên theo phân phối đã giả thiết. Chúng tôi cố định  $\mu=0.2$ ;  $t_{OFF \rightarrow MIDDLE} = 200s$  trong tất cả các thử nghiệm.

### 5.1 Thử nghiệm với ON/OFF/MIDDLE/ $\infty$ /STAG

Với thử nghiệm này, chúng tôi thực hiện để chứng minh hiệu quả của hệ thống có các máy chủ ở trạng thái MIDDLE với thuật toán chúng tôi đề xuất so với hệ thống chỉ có hai trạng thái ON, OFF. Hình 5.1a, 5.1b thể hiện số lượng máy chủ trung bình ở trạng thái MIDDLE và SETUP khi hệ thống được áp dụng thuật toán kiểm soát việc bật máy sang trạng thái MIDDLE. Số lượng máy chủ này tăng lên khi cường độ công việc đến hệ thống tăng lên, ở mức độ phù hợp nên có thể đảm bảo chi phí cho việc duy trì các máy chủ đó không quá lớn. Khi có máy chủ ở MIDDLE thời gian chờ đợi của công việc giảm đáng kể. Hình 5.3a, 5.3b thể hiện thời gian chờ đợi trung bình của các job khi  $\lambda$  thay đổi, tương ứng với  $\alpha=0.1$ ;  $0.01$ . Với  $\alpha=0.1$ , khi  $\lambda$  thay đổi từ 0.5 đến 10, thời gian chờ đợi trung bình của hệ thống không có máy chủ ở MIDDLE dao động trong khoảng từ 160s đến 200s lớn hơn nhiều so với hệ thống có trạng thái MIDDLE. Trường hợp  $\alpha=0.01$  tương tự, chênh lệch này vào khoảng 130s đến 190s xấp xỉ bằng khoảng thời gian để một máy bật từ OFF sang MIDDLE. Từ đó cho thấy giảm được thời gian trễ do bật máy chủ, chất lượng dịch vụ tăng lên.

### 5.2 Thử nghiệm với ON/OFF/MIDDLE/ $\infty$

Thời gian chờ đợi, số lượng máy chủ MIDDLE và SETUP trung bình được thể hiện trong hình 5.2 và hình 5.4. Số lượng máy chủ MIDDLE và SETUP tăng lên theo  $\lambda$  tương tự mô hình ON/OFF/MIDDLE/ $\infty$ /STAG. Trong mô hình này thời gian chờ đợi có xu hướng giảm khi cường độ công việc ( $\lambda$ ) tăng. Nguyên nhân do không hạn chế số lượng máy chủ có thể setup đồng thời nên khi số công việc tăng lên, lượng máy chủ SETUP cũng tăng lên, khả năng phục vụ của hệ thống tăng, thời gian chờ đợi trung bình giảm xuống. Các máy chủ được SETUP đồng thời nên lượng máy chủ MIDDLE và SETUP khá lớn so với mô hình không setup đồng thời (hình 5.2a và 5.2b). Hình 5.4 cho thấy thời gian chờ đợi của hệ thống có MIDDLE luôn nằm dưới hệ thống không có MIDDLE, rõ ràng khi được điều chỉnh bật máy chủ hợp lý khả năng phục vụ của hệ thống tăng lên.

## 6. KẾT LUẬN

Trong nghiên cứu này, chúng tôi đã đề xuất hệ thống với mô hình định nghĩa thêm trạng thái trung gian MIDDLE, xét với trường hợp tại bất kỳ thời điểm, chỉ có nhiều nhất một máy chủ trong chế độ SETUP và trường hợp không hạn chế số lượng máy chủ có thể setup đồng thời. Bằng cách sử dụng phân phối Poisson và chuỗi Markov, chúng tôi đã đưa ra thuật toán kiểm soát bật, tắt các máy OFF sang MIDDLE, thuật toán này nhằm đảm bảo luôn có máy chủ ở MIDDLE khi hệ thống cần và điều khiển số lượng máy ở MIDDLE để tránh lãng phí điện năng. Khi so sánh kết quả thực nghiệm giữa mô hình hệ thống có và không

có trạng thái MIDDLE thu được thông qua thí nghiệm mô phỏng bằng CloudSim [8], chúng tôi đã chứng minh hiệu quả đạt được trong việc giảm thời gian chờ đợi trung bình (mean waiting time) của hệ thống có trạng thái MIDDLE và tính toán được số lượng máy MIDDLE trung bình cần để đáp ứng yêu cầu giảm thời gian trễ do bật máy chủ. Tuy nhiên, chúng tôi mới chỉ dừng lại ở vấn đề tối thiểu thời gian chờ đợi trung bình. Đối với bài toán thực tế, chúng tôi còn phải giải quyết vấn đề cân bằng giữa tối thiểu thời gian chờ đợi trung bình và điện năng tiêu thụ. Trong tương lai, chúng tôi sẽ đi sâu hơn nữa để giải quyết vấn đề trên trong môi trường mô phỏng. Với kết quả có được, chúng tôi hi vọng nghiên cứu này sẽ đem lại một kết quả khả quan để có thể trở thành tiền đề nghiên cứu các chiến lược điều khiển và quản lý máy chủ tốt hơn, ứng dụng trong việc tăng chất lượng các dịch vụ điện toán đám mây.

## 7. LỜI CẢM ƠN

Trong công trình nghiên cứu này, chúng em xin cảm ơn sự hướng dẫn tận tình của TS. Nguyễn Bình Minh, cũng như sự giúp đỡ từ tập thể các thầy cô giáo ở viện Công nghệ thông tin và truyền thông. Chúng em cũng xin cảm ơn ban tổ chức cuộc thi SVNCKH đã giúp chúng em có một môi trường để làm quen với công việc nghiên cứu khoa học.

## 8. TÀI LIỆU THAM KHẢO

- [1] <http://www.gartner.com/newsroom/id/2867917>.
- [2] Quan, Dang Minh, et al. "Energy efficient resource allocation strategy for cloud data centres." *Computer and Information Sciences II*. Springer London, 2012. 133-141.
- [3] Buyya, Rajkumar, Anton Beloglazov, and Jemal Abawajy. "Energy-efficient management of data center resources for cloud computing: A vision, architectural elements, and open challenges." *arXiv preprint arXiv:1006.0308* (2010).
- [4] Long, Saiqin, Yuelong Zhao, and Wei Chen. "A three-phase energy-saving strategy for cloud storage systems." *Journal of Systems and Software* 87 (2014): 38-47.
- [5] Baliga, Jayant, et al. "Green cloud computing: Balancing energy in processing, storage, and transport." *Proceedings of the IEEE* 99.1 (2011): 149-167.
- [6] Beloglazov, Anton, Jemal Abawajy, and Rajkumar Buyya. "Energy-aware resource allocation heuristics for efficient management of data centers for cloud computing." *Future generation computer systems* 28.5 (2012): 755-768.
- [7] Phung-Duc, Tuan. "Server farms with batch arrival and staggered setup." *Proceedings of the Fifth Symposium on Information and Communication Technology*. ACM, 2014.
- [8] Calheiros, Rodrigo N., et al. "CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms." *Software: Practice and Experience* 41.1 (2011): 23-50.
- [9] Gandhi, A, Harchol-Balter, M. and Adan, I. (2010). "Server farm with setup cost", *Performance Evaluation*, 67, 1123-1138.
- [10] Gandhi, A, Harchol-Balter, M. and Adan, I. (2010). "Decomposition results for an M/M/k with staggered setup". *ACM SIGMETRICS Performance Evaluation Review*, 38, 48-50