

# QLA: روشی جدید برای ایجاد همکاری در سیستمهای چند عامله

برنا جعفر پور      محمدرضا میبدی

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

دانشگاه صنعتی امیر کبیر

تهران ایران

(jafarpour@aut.ac.ir, mmeybodi@aut.ac.ir)

## چکیده

اتوماتای یادگیر یک مدل انتزاعی است که تعداد محدودی عمل دارد که می‌تواند آنها را در محیط خود انجام دهد. اتوماتای یادگیر با توجه به پاسخی که از محیط در مقابل عمل خود میگیرد رفتار خود را به گونه ای تغییر می‌دهد که بیشترین پاداش را از محیط دریافت کند. در این مقاله یک مدل جدید بر پایه اتوماتاهای یادگیر برای ایجاد همکاری در سیستمهای چند عاملی پیشنهاد میگردد. در این مدل در هر حالت محیط یک اتوماتای یادگیر قرار داده می‌شود که حالت بعد را برای عاملی که در آن حالت قرار دارد تعیین می‌کند. برای بررسی کارایی این روش، از آن برای حل مسئله جمع آوری اشیاء استفاده شده است. آزمایشها نشان می‌دهد که این روش از کارایی بالاتری در مقایسه با روشهای برپایه ی Q-Learning و فرومون برخوردار می‌باشد.

**کلمات کلیدی:** سیستمهای چند عاملی، هماهنگی، اتوماتاهای یادگیر، یادگیری Q، جمع آوری اشیاء

## QLA: A New Coordination Method for Multi Agent Systems

B. Jafarpour      M. R. Meybodi

Computer Engineering and Information Technology Department

Amirkabir University of Technology

Tehran Iran

(jafarpour@aut.ac.ir, mmeybodi@aut.ac.ir)

## Abstract

A Learning automaton is an abstract model which randomly selects one action out of its finite actions and performs it on a random environment. Based on the responses received from the automaton changes its behavior so that it receives maximum reward from the environment. In this paper a new model based on learning automata for multi agent coordination is proposed. In this model, a learning automaton is put in each state of the environment which is responsible for determining the next state for the visiting agent. To show the performance of proposed model, it is used to solve multi agent Foraging Problem and compared the results with Q-Learning and Pheromone based solutions. Comparisons show the superiority of QLA over Q-Learning and Pheromones.

**Keywords:** Multi Agent Systems, Coordination, Learning Automata, Q-Learning, Foraging Problem.

## ۱- مقدمه

مختلفی از جمله الگوریتمهای تکاملی<sup>۱</sup>، اتوماتاهای یادگیر<sup>۲</sup> و فرومون<sup>۳</sup> استفاده شده است.

از پردازش تکاملی برای یافتن پارامترهای بهینه برای عاملها استفاده شده است [1-4]. در [1] از PSO<sup>۴</sup> برای یادگیری وزنیهای شبکه ی عصبی ای که رفتار رباتها را کنترل می‌کنند استفاده شده

ایجاد هماهنگی در یک سیستم چند عاملی امری بسیار مهم و در عین حال دشوار می‌باشد. چنانچه عاملها در یک سیستم همکاری لازم با یکدیگر را نداشته باشند، مزیتهایی که سیستمهای چند عاملی ارائه می کنند قابل دسترسی نخواهد بود. ایجاد هماهنگی بین عاملها عبارت است از همکاری عاملها به شکلی که پاداش دریافتی تیم در تعامل با محیط به صورت کلی (و نه صرفا به صورت فردی) افزایش پیدا کند. برای ایجاد هماهنگی بین عاملها در یک سیستم چند عاملی از روشهای

<sup>1</sup> Evolutionary Computation

<sup>2</sup> Learning Automata

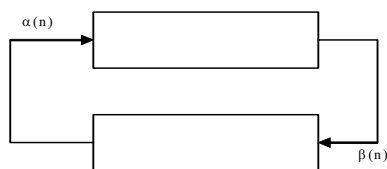
<sup>3</sup> Pheromone

<sup>4</sup> Particle Swarm Optimization

می‌گیرد در غیر این صورت از مفهوم آنتروپی برای تعیین پاداش و جریمه استفاده می‌شود. برای ارزیابی کارایی روش پیشنهادی از آن برای حل مسئله ی جمع آوری اشیاء<sup>11</sup> استفاده شده است. مقایسه ها با روشهای بر پایه ی فرومون و Q-Learning بعنوان یک روش یادگیری تقویتی نشان می‌دهد که روش پیشنهادی ابزار مناسبی برای ایجاد هماهنگی در سیستمهای چند عاملی می باشد. ادامه این گزارش بدین صورت سازماندهی شده است. در بخش ۲ اتوماتاهای یادگیر بطور خلاصه معرفی می‌شود. بخش ۳ روش پیشنهادی را مطرح می‌کند. در بخش ۴ مسئله ی جمع آوری اشیاء و روشهای حل آن برپایه فرومون و Q-Learning و روش پیشنهادی معرفی می‌شوند. نتایج شبیه سازی و نتیجه‌گیری به ترتیب در بخشهای ۵ و ۶ آمده است.

## ۲- اتوماتاهای یادگیر

اتوماتای یادگیر یک مدل انتزاعی است که تعداد معدودی عمل را می تواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی شده و پاسخی به اتوماتای یادگیر داده می شود. اتوماتای یادگیر از این پاسخ استفاده نموده و حالت درونی خود را به روز می کند و دوباره عمل خود را برای مرحله بعد انتخاب می کند [13]. اتوماتا با تعامل با محیط عمل بهینه را فراگیری می کند و به این شکل پاداش دریافتی خود را از محیط حد اکثر می کند. شکل ۱ ارتباط بین اتوماتای یادگیر و محیط را نشان می دهد.



شکل ۱- ارتباط بین اتوماتای یادگیر و محیط

محیط را می توان توسط یک سه تایی  $E \equiv \{\alpha, \beta, c\}$  نشان داد که در آن  $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه ورودیه ها،  $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\}$  مجموعه خروجیها و  $c \equiv \{c_1, c_2, \dots, c_r\}$  مجموعه احتمالهای جریمه  $\alpha$  ها می باشد. هر گاه  $\beta$  مجموعه دو عضوی باشد، محیط از نوع P می باشد. در چنین محیطی  $\beta_1 = 1$  به عنوان جریمه و  $\beta_2 = 0$  به عنوان پاداش در نظر گرفته می شود. در محیط از نوع Q،  $\beta(n)$  می تواند به طور گسسته یک مقدار از مقادیر محدود در فاصله  $[0,1]$  و در محیط از نوع S،  $\beta$  متغیر تصادفی در فاصله  $[0,1]$  است.  $c_i$  احتمال اینکه عمل  $\alpha_i$  نتیجه نامطلوب<sup>12</sup> داشته باشد، می باشد. در محیط ایستا<sup>13</sup> مقادیر  $c_i$  بدون تغییر می مانند، حال آن که در محیط غیر ایستا<sup>14</sup> این مقادیر در طی زمان تغییر می کنند.

است. در [2] از برنامه ریزی ژنتیک<sup>5</sup> برای تکامل یک تیم فوتبال استفاده شده است. الگوریتم ژنتیک در [3] برای تعیین پارامترهای مربوط به قرار دادن فرومون در تیمی از ربات ها استفاده شده است. در [4] از PSO برای تعیین پارامترهای یک تیم نا همگن از عاملها که عمل خوشه بندی را انجام می دهند استفاده شده است. اتوماتاهای یادگیر اخیرا بعنوان مدلی برای ایجاد هماهنگی در سیستمهای چند عامله مورد توجه محققان قرار گرفته است [5-9]. در [5] از اتوماتای یادگیر برای ایجاد هماهنگی در بین عاملهایی که عمل خوشه بندی را انجام می دهند استفاده شده است. در این مدل هر عامل دارای یک اتوماتای یادگیر است که تعیین می کند داده ی حمل شده توسط عامل باید در کجای صفحه ای که عمل خوشه بندی انجام می پذیرد قرار داده شود. نشان داده شده است که عامل ها بدون داشتن ارتباط مستقیم با یکدیگر قادر به ایجاد هماهنگی در قرار دادن داده ها بر روی صفحه می باشند. در [6] از اتوماتای یادگیر برای یادگیری اعمال بهینه در عاملهای یک تیم فوتبال استفاده شده است. در [7] نشان دهنده شده است که از اتوماتای یادگیر سلسله مراتبی<sup>6</sup> می توان به عنوان یک یادگیر تقویتی در سیستمهای چند عاملی مورد استفاده قرار گیرد. در [8] از اتوماتای یادگیر برای تعیین پارامتر بهینه کولونی مورچه ها<sup>7</sup> استفاده شده است. در [9] از اتوماتای یادگیر برای هماهنگی در اشتراک از یک خط مشترک ارتباطی در شبکه های بیسیم تک کاره<sup>8</sup> استفاده شده است. در این روش هر عامل دارای یک اتوماتای یادگیر می باشد که با توجه به داده های نویزی تخمین می زند چه عاملی در هر لحظه ارسال پیام خواهد کرد.

فرومونها یکی دیگر از مدلهایی است که برای ایجاد هماهنگی در سیستمهای چند عامله مورد استفاده قرار گرفته است. [10-12]. در [10] فرومونها برای تعیین کوتاهترین مسیر بین لانه و غذا در یک سیستم چند عاملی استفاده شده اند. در [11] نشان داده شده است که فرومونها گزینه ی مناسبی برای حل مسئله ی گشت زنی<sup>9</sup> در سیستمهای چند عاملی می باشند. در [12] نوع جدیدی فرومون برای سیستمهای چند عاملی معرفی و در حل مسئله ی  $A^*$  فیزیکی<sup>10</sup> استفاده شده است.

در این مقاله یک مدل جدید بر پایه اتوماتاهای یادگیر برای ایجاد همکاری در سیستمهای چند عاملی پیشنهاد میگردد. در این مدل در هر حالت محیط یک اتوماتای یادگیر قرار داده می شود. وظیفه ی اتوماتای یادگیر تعیین حالت بعدی برای عاملی است که در آن حالت قرار دارد. چنانچه اتوماتای یادگیر عامل را به حالت هدف ببرد، پاداش

<sup>5</sup> Genetic Programming

<sup>6</sup> Hierarchical

<sup>7</sup> Ant Colony

<sup>8</sup> Ad hoc Wireless Network

<sup>9</sup> Patrolling Problem

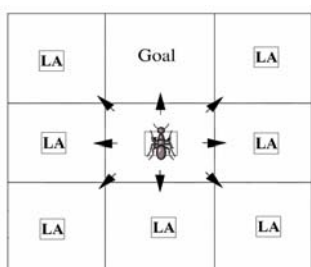
<sup>10</sup> Physical A\*

<sup>11</sup> Foraging

<sup>12</sup> Unfavorable

<sup>13</sup> Stationary

<sup>14</sup> Non-Stationary



شکل ۲- چینش اتوماتای یادگیر در حالت‌های محیط

در ابتدا اتوماتاهای یادگیر تمام عمل‌های خود را با احتمالی یکسان انتخاب می‌کنند. چنانچه عمل اتوماتا منجر به ورود عامل به حالت هدف شود اتوماتا پاداش می‌گیرد، در غیر اینصورت از آنتروپی بردار احتمال اتوماتای یادگیر حالت بعد برای تعیین پاداش یا جریمه استفاده می‌شود. آنتروپی بردار احتمال میزان عدم اتوماتای یادگیر حالت بعد را در انتخاب عمل خود نشان می‌دهد. هر چه آنتروپی بیشتر باشد میزان عدم قطعیت بیشتر است. عدم قطعیت بالا در بردار احتمال اتوماتای یادگیر به این معنی است که این اتوماتای یادگیر دارای اطلاعات مفیدی برای رسیدن به هدف نیست و عمل‌های خود را به صورت تصادفی انتخاب می‌کند. ولی چنانچه عدم قطعیت کم باشد به این معنی است که اتوماتای یادگیر با احتمال بالایی یکی از اعمال خود را انتخاب می‌کند و دارای اطلاعات مفیدی برای رسیدن به هدف می‌باشد. فرض کنید که  $\{p_1, p_2, \dots, p_r\}$  بردار احتمال اعمال یک اتوماتای یادگیر باشد. آنتروپی این بردار احتمال به شکل زیر تعیین می‌شود.

$$E(P) = -\sum_{i=1}^r p_i \log(p_i) \quad (2)$$

هرچه میزان عدم قطعیت بیشتر باشد مقدار آنتروپی بیشتر خواهد بود. زمانی  $E$  بیشترین مقدار را خواهد داشت که تمام اعمال احتمالی یکسان داشته باشند ( $P_{equal} = \{p_1=p_2=\dots=p_r=1/r\}$ ) و زمانی کمترین مقدار (برابر با ۰) را خواهد داشت که  $\exists i \ p_i = 1 \wedge \forall j \neq i \ p_j = 0$ . برای اینکه مقدار آنتروپی را به مقداری بین ۰ و ۱ تبدیل کنیم تا به عنوان بردار تقویتی در اتوماتای ساختار متغیر مدل S قابل استفاده باشد از فرمول زیر استفاده می‌شود.

$$\beta = E((P) / E(P_{equal}))^K \quad (3)$$

K پارامتر روش می‌باشد. میزان این پارامتر باید به دقت تعیین شود. مقادیر بالای این پارامتر باعث می‌شود که  $\beta$  ها مقادیر کمی داشته باشند و اتوماتاهای یادگیر بیش از حد پاداش ببینند. این امر باعث می‌شود که انتخاب یک عمل با احتمال بیشتری پاداش بگیرد. به این شکل اتوماتاها با احتمال بیشتری حالت‌های نه چندان مناسب را انتخاب

اتوماتای یادگیر به دو گروه با ساختار ثابت و با ساختار متغیر تقسیم می‌گردد [13]. اتوماتای یادگیر با ساختار متغیر<sup>۱۵</sup> توسط ۴ تایی  $\{\alpha, \beta, p, T\}$  نشان داده می‌شود که در آن  $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه عمل‌های اتوماتا،  $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\}$  مجموعه ورودی‌های اتوماتا و  $p = \{p_1, p_2, \dots, p_r\}$  بردار احتمال انتخاب هر یک از اعمال و  $p(n+1) = T[\alpha(n), \beta(n), p(n)]$  الگوریتم یادگیری می‌باشد. در این نوع از اتوماتاها، اگر عمل  $\alpha_i$  در مرحله n ام انتخاب شود و این عمل، پاسخ مطلوب از محیط دریافت نماید، احتمال  $p_i(n)$  افزایش یافته و سایر احتمالات کاهش می‌یابند. برای پاسخ نامطلوب احتمال  $p_i(n)$  کاهش یافته و سایر احتمالات افزایش می‌یابند. این تغییرات به گونه ای صورت می‌پذیرد که جمع احتمالات برابر با یک باقی بماند. فرمول (۴) یکی از الگوریتم‌های یادگیری خطی در اتوماتای با ساختار متغیر مدل S را نشان می‌دهد. فرض کنید که عمل i توسط اتوماتای یادگیر انتخاب شده باشد و مقدار  $\beta$  توسط محیط به اتوماتا داده شده باشد :

$$\begin{aligned} p_i(n+1) &= p_i(n) + a(1-\beta(n))(1-p_i(n)) - b\beta(n)p_i(n) \\ p_j(n+1) &= p_j(n) - a(1-\beta(n))p_j(n) + b\beta(n)\left(\frac{1}{r-1} - p_j(n)\right) \quad \forall j, j \neq i \end{aligned} \quad (1)$$

a و b به ترتیب پارامتر پاداش و پارامتر جریمه می‌باشد. زمانی که a و b با هم برابر باشند، الگوریتم  $L_{R-P}$ <sup>۱۶</sup>، زمانی که a از b خیلی کوچکتر باشد، الگوریتم  $L_{R-\epsilon P}$ <sup>۱۷</sup> و زمانی که b مساوی صفر باشد، الگوریتم  $L_{R-I}$ <sup>۱۸</sup> نامیده می‌شود.

### ۳- روش پیشنهادی

در این بخش روش پیشنهادی برای یادگیری در سیستم‌های چند عاملی به کمک اتوماتای یادگیر را بررسی می‌کنیم. در این مدل در هر حالت محیط یک اتوماتای یادگیر با ساختار متغیر مدل S قرار داده می‌شود. تعداد اعمال اتوماتا در هر حالت برابر است با تعداد حالت‌های همسایه ی آن حالت و هر عمل متناظر با یکی از حالات همسایه می‌باشد. هر عامل برای تعیین حالت بعدی خود از اتوماتای یادگیر کمک می‌گیرد؛ به این معنی که عامل به حالت متناظر با عمل انتخاب شده توسط اتوماتا می‌رود. شکل ۲ چگونگی چینش اتوماتای یادگیر در حالت‌های یک محیط نشان می‌دهد.

<sup>15</sup> Variable Structure Learning Automata

<sup>16</sup> Linear Reward Penalty

<sup>17</sup> Linear Reward Epsilon Penalty

<sup>18</sup> Linear Reward Inaction

#### ۴-۱ تعریف مسئله

در این مسئله تعدادی عامل بر روی صفحه ای دو بعدی مانند اتوماتای سلولی قرار دارند. عامل ها از یک لانه شروع به جستجو برای یافتن غذا (اشیاء) می کنند. هنگامی که این عاملها به منبع غذا برسند یک واحد از آن را برمی دارند و سعی می کنند آن را به لانه تحویل دهند. عاملها تنها در سلولهای صفحه می توانند حرکت کنند و بر خلاف اتوماتای سلولی صفحه مدور<sup>۲۰</sup> نمی باشد. یعنی عامل نمی تواند از یک طرف صفحه خارج شود و از طرف دیگر وارد شود. کارایی یک تیم از عامل ها در این مسئله برابر است با تعداد واحدهایی از غذا که در مدت زمان مشخصی به لانه تحویل داده می شود. بدیهی است که هرچه مسیری که عاملها از لانه به غذا و برعکس طی می کنند کوتاهتر باشد و عاملهای بیشتری در جمع آوری اشیاء شرکت کنند کارایی تیمی بیشتر خواهد شد. این مسئله را می توان به دو زیر مسئله ی یافتن کوتاهترین مسیر از لانه به غذا و از غذا به لانه تقسیم کرد. در این مسئله عاملها حافظه ای ندارند و قادر تشخیص مکان خود در صفحه نمی باشند و تنها به اطلاعاتی که در سلولهای همسایه توسط دیگر عاملها نوشته شده دسترسی دارند. هر عامل می تواند این داده ها را بخواند و با توجه به آنها برای مسیر حرکت خود تصمیم گیری کند. علاوه بر این، عامل می تواند این داده ها را تغییر دهد. عاملها باید به کمک هم و با نوشتن داده های مناسب با یکدیگر همکاری کرده تا کوتاهترین مسیر از لانه به غذا و برعکس را بیابند. برای اطلاعات بیشتر در مورد این مسئله می توانید به [14] مراجعه کنید. در این مقاله سه روش را برای حل این مسئله بررسی و آنها را با هم مقایسه می کنیم. این ۳ روش عبارتند از Q-Learning، فرومونها و روش پیشنهادی.

#### ۴-۲ حل مسئله با Q-Learning

Q-Learning [15] یکی از روشهای یادگیری تقویتی است که به علت کارایی بالا و سادگی استفاده در کاربردهای فراوانی مورد استفاده قرار گرفته است. در این روش در هر حالت مسئله یک آرایه به نام Q وجود دارد. طول این آرایه برابر است با تعداد حالتهای همسایه ی آن حالت. هر خانه ی آرایه متناظر است با یکی از حالتهای همسایه و مقدار موجود در آن خانه میزان مطلوبیت همسایه ی متناظر را نشان می دهد. فرض کنید که یک عامل در حالت S قرار داشته باشد. عامل برای انتخاب حالت بعد با توجه به مقادیر Q تصمیم گیری می کند. یکی از روشهای متداول برای انتخاب عمل برابر است با انتخاب حالت بعد با احتمال متناسب با مقدار Q متناظر با آن. به عنوان مثال به شکل زیر :

$$P(s, a) = Q(s, a)^K / \sum_{a \in A_s} Q(s, a)^K \quad (5)$$

می کنند و این به معنی جستجوی<sup>۱۹</sup> بیشتر در محیط است. زمانی که K مقدار بالایی داشته باشد این احتمال وجود دارد که عاملها تشکیل یک حلقه دهند و به دنبال یکدیگر شروع به حرکت کنند. به این شکل عاملها ممکن است در یک حلقه بسته قرار بگیرند. هر چه میزان K کمتر باشد  $\beta$  ها مقادیر بیشتری خواهند داشت. این امر باعث می شود که اتوماتاها بیشتر جریمه شوند. این مسئله باعث می شود که حتی حالتی مطلوب نیز پاداش لازم را نگیرند. به این شکل عاملها قادر به یافتن مسیر مناسبی به سمت هدف نخواهند بود و در صفحه بی هدف پرسه می زنند. برای تعیین میزان نرخ جستجو (Exploration) و بهره برداری (Exploitation) می توان از تغییر پارامتر یادگیری (a) در فرمول (۱) استفاده کرد. هرچه نرخ یادگیری بیشتر باشد عاملها حالتی خوب را با احتمال بیشتری در ادامه انتخاب خواهند کرد ولی چنانچه مقدار این پارامتر کم باشد، عاملها دیرتر به سمت حالتی خوب جذب می شوند و بیشتر به جستجو می پردازند. فرض کنید که عامل در حالت S باشد و اتوماتای یادگیر آن ( $LA_s$ ) عامل را به حالت  $S'$  هدایت کند. در اینصورت تعیین سیگنال تقویتی به این شکل تعیین می شود :

$$\beta_s = \begin{cases} 1 & S' = \text{Out of bound} \\ 0 & S' = \text{Goal} \\ (E(P(LA_{S'})) / E(P_{equal}))^K & \text{otherwise} \end{cases} \quad (4)$$

$P(LA_s)$  نشان دهنده ی بردار احتمال اتوماتای یادگیر قرار گرفته در حالت S می باشد. چنانچه چندین عامل در یک حالت قرار داشته باشند عاملها به صورت تصادفی انتخاب می شوند و انتخاب حالت بعدی و آموزش اتوماتای یادگیر انجام می شود. به علت شباهتی که این روش به Q-Learning دارد آن را QLA نام گذاری کرده ایم. شبه کد این الگوریتم را می توان به شکل زیر نوشت :

```

Initialize.
While not done do
  For each agent i do in parallel
    Activate  $LA_i$  (LA residing in current state of agent i ( $S_i$ ))
    Move to  $S'_i$  (next state selected by  $LA_i$ )
    Compute  $\beta_i$  (reinforcement signal of  $LA_i$ )
    Train LA residing in  $S_i$  according to  $\beta_i$ 
  End parallel for
End while

```

شکل ۳- شبه کد QLA

#### ۴-۳ مسئله ی جمع آوری اشیاء

در این بخش مسئله ی جمع آوری اشیاء را مطرح می کنیم و توضیح می دهیم که این مسئله چگونه توسط Q-Learning، فرومون ها و روش پیشنهادی قابل حل می باشد.

<sup>20</sup> Wraparound

<sup>19</sup> Exploration

در این الگوریتم هرچه میزان  $K$  بیشتر باشد الگوریتم بیشتر به بهره برداری (Exploitation) و هرچه کمتر باشد بیشتر به جستجو (Exploration) می پردازد. فرض کنید که عامل در حالت  $s$  باشد و بردار  $Q$  عامل را به حالت  $s'$  هدایت کند. در اینصورت بردار  $Q$  حالت  $s$  به این شکل به روز می شود :

$$Q(s, a) = r + \gamma \max_{a' \in A_{s'}} (Q(s', a')) \quad (6)$$

$A_s$  برابر است با مجموعه اعمال ممکن در حالت  $s$  و  $\gamma$  پارامتر الگوریتم می باشد که معمولاً مقداری نزدیک به یک در نظر گرفته می شود. برای حل مسئله ی جمع آوری اشیاء توسط این روش در هر سلول صفحه دو آرایه ی  $Q$  با طول ۸ قرار داده می شود.  $Q_1$  برای یافتن کوتاهترین مسیر از لانه به سمت غذا و  $Q_2$  برای مسیر برعکس استفاده می شود. چنانچه عامل حامل غذا نباشد به سلول غذا وارد شود و یا چنانچه عامل حامل غذا باشد و به سلول لانه وارد شود پاداش می گیرد. توجه کنید که زمانی که عاملها حامل غذا باشند تنها به  $Q_2$  دسترسی دارد و آن را به روز می کند در غیر اینصورت تنها به  $Q_1$  دسترسی دارد و آن را به روز می کند.

### ۳-۴ حل مسئله با فرومون

برای حل این مسئله از فرومونها نیز استفاده شده است [10]. در این روش عاملها زمانی که حامل غذا نباشند فرومون "به سمت لانه" بر روی سلولها قرار می دهند و فرومون "به سمت فرومونهای غذا" را دنبال می کنند. عاملها هر چه از لانه دور شوند میزان فرومونی که بر روی سلولها قرار می دهند کاهش می یابد. زمانی که عامل به منبع غذا برسد یک واحد از آن را بر می دارد و پس از آن در هنگام حرکت، فرومون "به سمت غذا" بر روی سلولها قرار می دهد و فرومون "به سمت لانه" را دنبال می کند. لانه و غذا همیشه به ترتیب دارای بیشترین مقدار فرومون به سمت لانه و فرومون به سمت غذا می باشند. عاملها در این روش برای قرار دادن فرومون از فرمول زیر استفاده می کنند :

$$Phr(s) = \max_{s' \in N(s)} (Phr(s')) - 2 \quad (7)$$

$Phr(s)$  برابرست با مقدار فرومون در سلول  $s$  و  $N(s)$  برابر است با ۸ همسایه مجاور سلول  $s$ . عاملها سلولهای همسایه را با احتمالی متناسب با مقدار فرومون در آن سلولها از طریق فرمول زیر انتخاب می کنند.

$$p(s') = (phr(s') + c)^K \quad (8)$$

همانند Q-Learning در این الگوریتم هرچه میزان  $K$  بیشتر باشد الگوریتم بیشتر به بهره برداری (Exploitation) و هرچه کمتر باشد

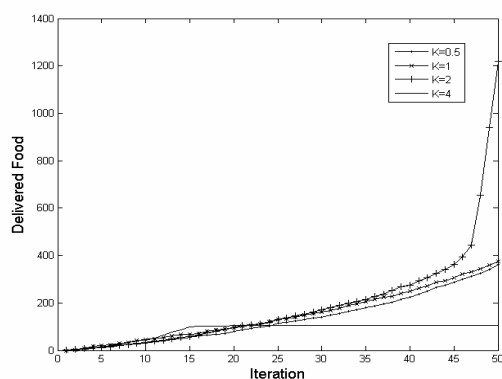
بیشتر به جستجو (Exploration) می پردازد. برای جلوگیری از گیر افتادن در مینیمم محلی در هر تکرار فرومونها در ضربی به نام  $\rho$  که کوچکتر از ۱ می باشد ضرب می شوند.

### ۳-۴ حل مسئله با QLA

برای حل این مسئله در هر سلول صفحه دو اتوماتای یادگیر با ۸ عمل قرار می دهیم. هر عمل متناظر است با یکی از خانه های همسایه. یک اتوماتا وظیفه هدایت عاملها به سمت غذا و دیگری وظیفه هدایت عاملها به سمت لانه را دارد. بسته به اینکه عامل حامل غذا باشد یا خیر، یکی از دو اتوماتای یادگیر غذا یا لانه وظیفه ی هدایت عامل را بر عهده دارد. همانند شکل ۲ اتوماتای یادگیر در سلولهای صفحه قرار می گیرند و زمانی که یک عامل در آنها قرار بگیرد اتوماتای یادگیر محل بعدی آن را در صفحه تعیین می کند. عامل نیز با توجه به به فرمول (۴) مقدار سیگنال تقویتی را تعیین می کند و با توجه به فرمول (۱) بردار احتمال اتوماتای یادگیر را به روز می کند.

### ۵- نتایج شبیه سازیها

در آزمایش اول تاثیر پارامتر  $K$  را بر روی الگوریتم پیشنهادی بررسی کرده ایم. در این آزمایش از یک صفحه ی  $30 \times 30$  استفاده کرده ایم. لانه در خانه ی (۵،۵) و غذا در خانه ی (۲۵،۲۵) قرار دارد. در QLA از روش یادگیری  $LRI$  با  $a$  و  $b$  به ترتیب برابر با ۰.۳ و ۰ استفاده شده است.  $K$  در تکرارهای مختلف برابر با ۰.۵، ۱، ۲ و ۴ قرار داده شده است. تعداد عاملها در تمام آزمایشها برابر با ۱۰ قرار داده شده است. شکل ۴ میزان غذای تحویل داده شده به لانه را در ۵۰۰۰ تکرار می بینید. شکل ۴ میانگین ۲۰ تکرار را گزارش می کند.

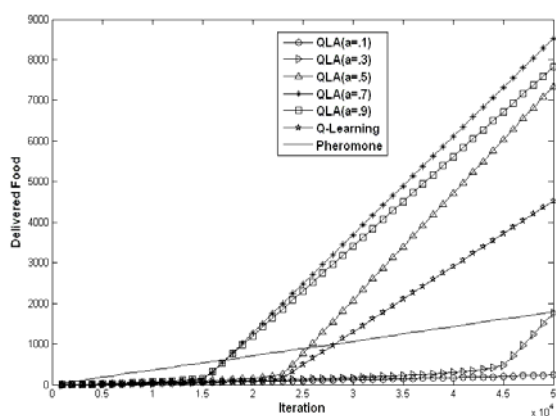


شکل ۴ - میانگین غذای تحویل داده شده به لانه توسط QLA با K های مختلف

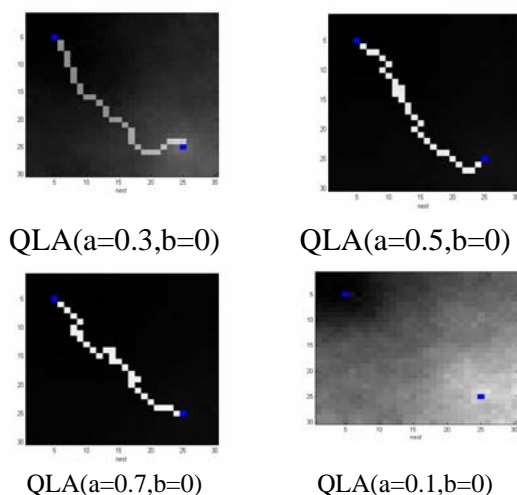
همانطور که می بینید  $K=2$  در این آزمایش بهترین جواب را تولید کرده است. بقیه ی مقادیر  $K$  برای الگوریتم مناسب نمی باشند. برای نمایش مسیر حرکت عامل ها بر روی آنها را به فرومون مجهز کرده ایم.

می‌شود که در عاملها در یک حلقه به دنبال یکدیگر شروع به حرکت کنند(شکل ۶ پایین سمت راست).

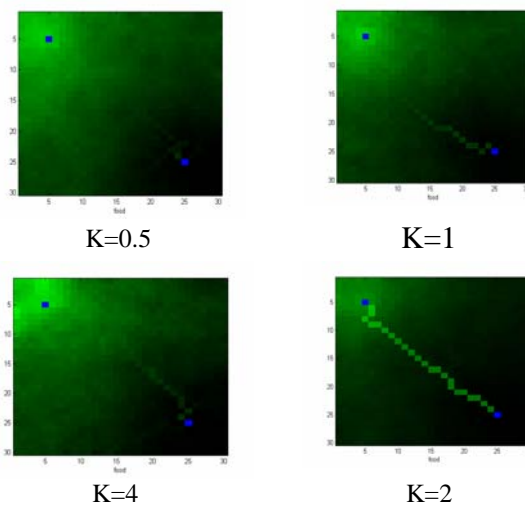
برای بررسی بیشتر ۲ آزمایش دیگر ترتیب داده ایم. در آزمایش اول عاملها در یک صفحه ی  $30 \times 30$  می‌توند حرکت کنند. لانه در خانه (۵،۵) و غذا در خانه ی (۲۵،۲۵) قرار دارد. مقدار  $K$  در روش  $Q$ -Learning، فرومون و  $QLA$  به ترتیب برابر با ۱۰، ۱۰ و ۲ قرار داده شده است. در روش بر پایه ی فرومون مقدار  $c$  برابر با 0.001 قرار داده شده است. در روش  $QLA$  از روش یادگیری  $LRI$  با  $a$  های برابر با 0.1، 0.3، 0.5، 0.7 و 0.9 و  $b$  برابر با 0 استفاده شده است تا تاثیر نرخ یادگیری در عملکرد الگوریتم را ببایم. در روش  $Q$ -Learning میزان پاداش ( $r$ ) برای ورود به لانه و غذا برابر ۵ و  $\gamma$  برابر با 0.95 قرار داده شده است. در شکل ۷ میزان غذای تحویل داده شده به لانه را در ۵۰۰۰ تکرار می‌بینید. شکل ۷ میانگین ۲۰ تکرار را گزارش می‌کند.



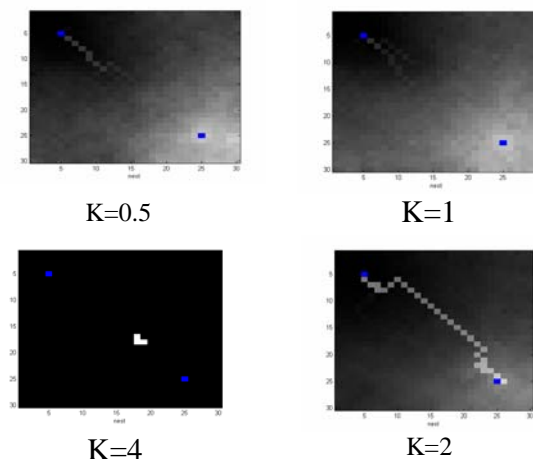
شکل ۷- میانگین غذای تحویل داده شده به لانه توسط سه روش مورد بررسی در صفحه ی  $30 \times 30$



دقت کنید که این فرومونها در  $QLA$  صرفا جنبه ی نمایش دارند و هیچ تاثیری در الگوریتم ندارند. عاملها زمانی که حرکت می‌کنند بر روی سلولها فرومون قرار می‌دهند. به این شکل روی مسیری که عاملها طی می‌کنند فرومون بیشتری جمع می‌شود. با تبدیل این فرومونها به تصویر، مسیری که عاملها در هر روش طی می‌کنند مشخص می‌شود. شکل ۵ مسیر طی شده از غذا به سمت لانه و شکل ۶ مسیر طی شده از لانه به سمت غذا به ازای  $K$  های مختلف را در انتهای ۵۰۰۰ اجرا نشان می‌دهد. در شکلهای ۵ و ۶ هرچه میزان فرومون در یک سلول بیشتر باشد آن سلول روشنتر نشان داده شده است.



شکل ۵- مقایسه ی مسیر طی شده از لانه به سمت غذا توسط  $QLA(a=3, b=0)$  با  $K$  های مختلف

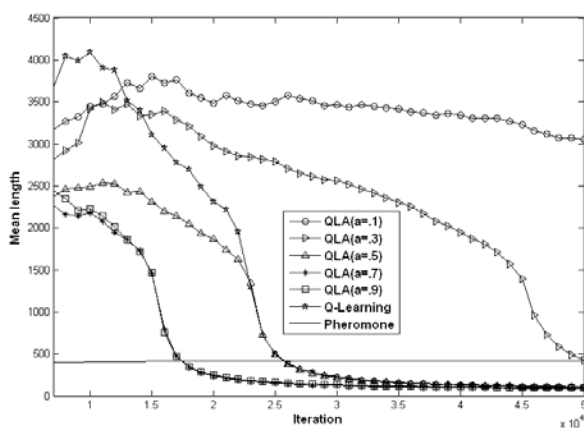


شکل ۶- مقایسه ی مسیر طی شده از غذا به سمت لانه توسط  $QLA(a=3, b=0)$  با  $K$  های مختلف

همانطور که در شکلهای ۵ و ۶ می‌بینید  $K=2$  در این آزمایش بهترین مسیر را یافته است. این شکلهای نتیجه ی شکل ۴ را تأیید می‌کنند. همانطور که در بخش ۳ مطرح کردیم مقادیر بالای  $K$  باعث

دارند و مسیرهای خوبی می‌یابند. البته نرخ بالای یادگیری باعث گرفتار شدن در مینیمم محلی نیز می‌شود. همانطور که در شکل بالا می‌بینید نرخ یادگیری 0.9 با وجود همگرایی سریع، نسبت به مقدار 0.7 دارای کارایی کمتری می‌باشد. دلیل این امر را می‌توان گرفتار شدن در مینیمم محلی دانست. راز کارایی بالای QLA در مسیری است که عاملها توسط این روش بین لانه و غذا می‌پیمایند. همانند آزمایش قبل از فرومونها برای نمایش مسیر طی شده توسط عاملها استفاده کرده ایم. شکل ۸ مسیر طی شده از غذا به سمت لانه و شکل ۹ مسیر طی شده از لانه به سمت غذا در ۳ روش مورد بررسی را در انتهای ۵۰۰۰ اجرا نشان می‌دهد.

همانطور که مشاهده میشود نرخهای یادگیری بزرگتر از 0.3 قادر به یافتن مسیرهای خوبی بین لانه و غذا می‌باشند. البته همانطور که در قبل مطرح کردیم مسیرهای یافته شده نشان می‌دهد که نرخ یادگیری 0.9 با وجود همگرایی بالا گاهی در مینیمم محلی گرفتار می‌شود. مسیر طی شده توسط QLA به خط راست بین لانه و غذا نزدیک است. عاملها تقریباً هیچگاه از مسیر پیدا شده خارج نمی‌شوند. شکل‌های ۸ و ۹ نشان می‌دهد روش پر پایه ی فرومون ها و Q-Learning جواب مناسبی را نتوانسته اند بیابند. برای مقایسه ی بیشتر، در تمام روش ها میانگین مسیر طی شده برای آوردن غذا ها به لانه در هر تکرار محاسبه شده است. شکل ۱۰ میانگین این پارامتر را در ۲۰ تکرار نشان می‌دهد.



شکل ۱۰ - میانگین مسیر طی شده برای تحول غذا به لانه توسط سه روش مورد بررسی در طول اجرا در صفحه ی ۳۰×۳۰

شکل ۱۰ نشان می‌دهد که روش پر پایه ی فرومونها به سرعت مسیر نسبتاً خوبی بین لانه و غذا پیدا می‌کند ولی پس از آن قادر به بهبود مسیر یافته شده نیست. پس از آن سریعترین همگرایی مربوط به QLA با نرخ های یادگیری 0.7 و 0.9 می‌باشد.

آزمایش قبل را برای صفحه ای ۶۰×۶۰ نیز تکرار کرده ایم. لانه در خانه (۱۵،۱۵) و غذا در خانه ی (۴۵،۴۵) قرار دارد. از QLA(a=0.7,b=0) استفاده کرده ایم. سایر پارامترها مشابه آزمایش



Phormone

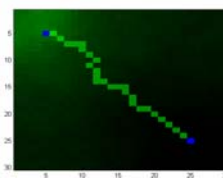


QLA(a=0.9,b=0)

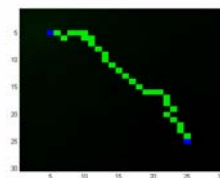


Q-Learning

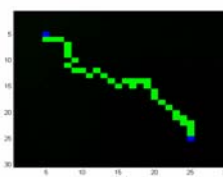
شکل ۸- مقایسه ی مسیر طی شده از غذا به سمت لانه در ۳ روش مورد بررسی در صفحه ی ۳۰×۳۰



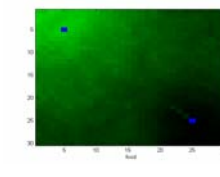
QLA(a=0.3,b=0)



QLA(a=0.5,b=0)



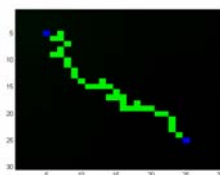
QLA(a=0.7,b=0)



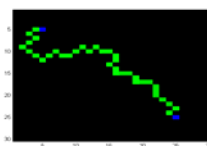
QLA(a=0.1,b=0)



Phormone



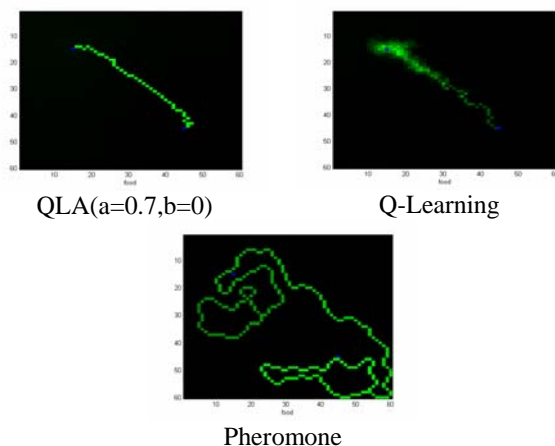
QLA(a=0.9,b=0)



Q-Learning

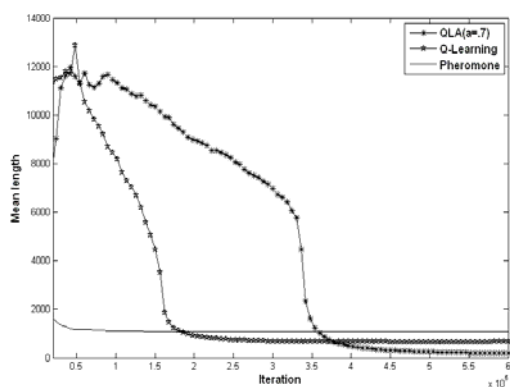
شکل ۹- مقایسه ی مسیر طی شده از لانه به سمت غذا در ۳ روش مورد بررسی در صفحه ی ۳۰×۳۰

همانطور که مشاهده میشود روش ارائه شده با نرخ یادگیری بالا دارای کارایی بسیار بالایی در مقایسه با Q-Learning و فرومونها می‌باشد. از میان نرخهای یادگیری، مقدار 0.7 و 0.9 همگرایی سریعی



شکل ۱۳ - مقایسه ی مسیر طی شده از لانه به سمت غذا در ۳ روش مورد بررسی در صفحه ی ۶۰×۶۰

همانطور که مشاهده میشود زمانی که طول مسیر بین لانه و غذا زیاد شده است افت مقادیر Q باعث سرگردانی عاملها پس از خروج از لانه و پرسه زدن آنها می شود. برای مقایسه ی بیشتر، در تمام روش ها میانگین مسیر طی شده برای آوردن غذا ها به لانه در هر تکرار محاسبه شده است. شکل ۱۴ میانگین این پارامتر را در ۲۰ تکرار نشان می دهد.



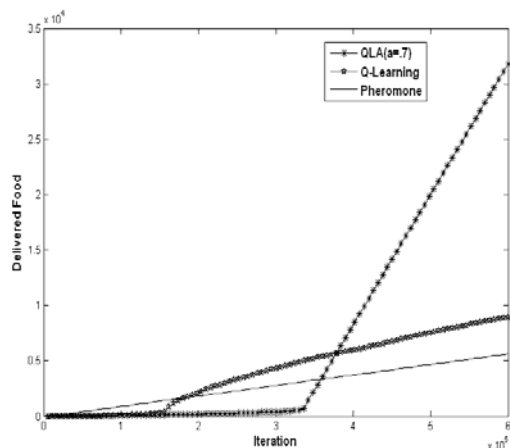
شکل ۱۴ - میانگین مسیر طی شده برای تحول غذا به لانه توسط سه روش مورد بررسی در طول اجرا در صفحه ی ۶۰×۶۰

همانند آزمایش قبل مشاهده میشود که QLA مسیر بهتری نسبت به Q-Learning و فرومون ها می یابد و غذای بیشتری جمع آوری می کند.

## ۶- نتیجه گیری

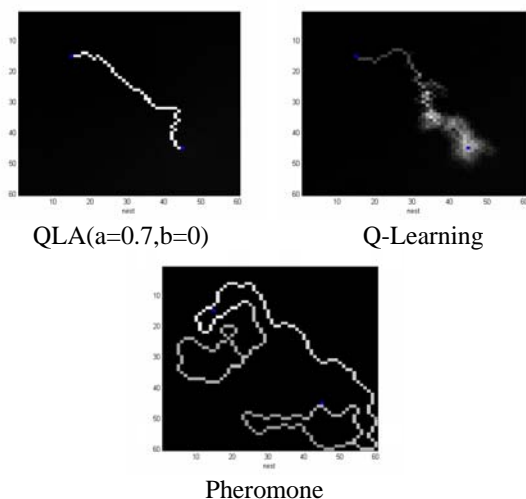
در این مقاله روشی برای ایجاد هماهنگی در بین عاملهای یک سیستم چند عاملی بر پایه ی اتوماتای یادگیر پیشنهاد و سپس از آن برای حل مسئله ی جمع آوری اشیاء استفاده گردید. مقایسه روش پیشنهادی با روشهای Q-Learning و فرومونها نشان داد که روش

قبل استفاده شده اند. در شکل ۱۱ میزان غذای تحویل داده شده به لانه را در ۶۰۰۰۰۰ تکرار می بینید. شکل ۱۱ میانگین ۲۰ تکرار را گزارش می کند.



شکل ۱۱ - میانگین غذای تحویل داده شده به لانه توسط سه روش مورد بررسی در صفحه ی ۶۰×۶۰

همانطور که مشاهده میشود  $K=2$ ،  $a=0.7$  و  $b=0.0$  پارامترهای مناسبی برای روش ارائه شده می باشند شکل ۱۲ و ۱۳ به ترتیب نمونه مسیر یافته شده از غذا به لانه و برعکس را در ۳ روش مورد بررسی نشان می دهد.



شکل ۱۲ - مقایسه ی مسیر طی شده از غذا به سمت لانه در ۳ روش مورد بررسی در صفحه ی ۶۰×۶۰



[15] Watkins, C., "Learning from Delayed Rewards", Thesis, University of Cambridge, England, 1989.

پیشنهادی ابزار مناسبی برای ایجاد هماهنگی در سیستمهای چند عاملی می باشد. دلیل کارایی بالای روش پیشنهادی را می توان در استفاده از مزیت هر دو روش Q-Learning و فرمونها دانست.

## مراجع

- [1] Pugh, J. and Martinoli, A., "Multi-Robot Learning with Particle Swarm Optimization", Proceedings of International Conference on Autonomous Agents and Multiagent Systems, pp. 441 - 448, 2006.
- [2] Haynes, T and Sen, S., "Crossover operators for evolving a team". In J. R. Koza, K. Deb, M. Dorigo, D. B. Fogel, M. Garzon, H. Iba and R. L. Riolo, editors, Genetic Programming 1997: Proceedings of the Second Annual Conference, pages 162–167, Stanford University, CA, USA, 13-16 July 1997.
- [3] Sauter, J., Matthews, R., Van Dyke Parunak, H. and Brueckner, S., "Evolving adaptive pheromone path planning mechanism"s. In Proceedings of First International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS-02), pp. 434-440, 2002.
- [4] Jafarpour, B. and Meybodi, M. R., "Adaptation of Ant Clustering Parameters using CLA-PSO ", The 1st National Data Mining Conference (IDMC'2007), Tehran, Iran, 2007.
- [5] Khojasteh, M. R. and Meybodi, M. R., "Evaluating Learning Automata as a Model for Cooperation in Complex Multi-Agent Domains", Lecture Notes in Artificial Intelligence, Springer Verlag, LNAI 4434, pp. 409-416, 2007.
- [6] Jafarpour, B. and Meybodi, M. R., "Improving Ant based Clustering Using Learning Automata", The 13<sup>th</sup> International Computer Society of Iran Conference (CSICC'2008), Kish Island, Iran, 2008.
- [7] Ann Nowé, Katja Verbeeck and Maarten Peeters: "Learning Automata as a Basis for Multi Agent Reinforcement Learning", LAMAS 2005.
- [8] Abdali, F. and Meybodi, M. R., "Adaptation of Ants Colony Parameters Using Learning Automata", Proceedings of 10th Annual CSI Computer Conference Iran, Telecommunication Research Center, Tehran, Iran, pp. 972-980, Feb. 2005.
- [9] Aggarwal, P. S. and Chenhui Liu, "Coordination within Multiple Learning Automata Agents: a Novel Distributed Permission Switching Protocol", International Conference on Integration of Knowledge Intensive Multi-Agent Systems, 2005.
- [10] Panait, L. and Luke, S., "Ant Foraging Revisited", Proceedings of the Ninth International Conference on the Simulation and Synthesis of Living Systems (ALIFE9). 2003.
- [11] Chu, Hoang-Nam, Glad, A., Simonin, O., Sempé, F., Drogoul, A. and Charpillet, F., "Swarm Approaches for the Patrolling Problem, Information Propagation vs. Pheromone Evaporation", Proceedings of the 19<sup>th</sup> IEEE International Conference on Tools with Artificial Intelligence (ICTAI'2007), Patras, Greece, October, 2007.
- [12] Felner, A., Shoshani, Y., Wagner, I. A. and Bruckstein, A. M., "Multi-agent Physical A\* using Large Pheromones", Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems, 2004. AAMAS 2004.
- [13] Narendra, K. S. and Thathachar, M. A. L., Learning Automata: An introduction, Prentice Hall, 1989.
- [14] Sen, S and Weiss, G., Multi-Agent Systems: A Modern Approach to Distributed Artificial Intelligence, MIT Press, 1999.