

# بهبود مدل مفهومی CLAAaS برای تحلیل داده‌های عظیم در ابر

لیلا محمدی نافچی<sup>۱\*</sup>، محمدعلی منتظری<sup>۲</sup>، علی آهون‌منش<sup>۳</sup>

<sup>۱</sup> دانشجوی مقطع کارشناسی ارشد، مؤسسه آموزش عالی صنعتی فولاد، فولادشهر،

Lm.eng@live.com

<sup>۲</sup> استادیار گروه مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی اصفهان،

montazeri@cc.iut.ac.ir

<sup>۳</sup> استاد، رئیس مؤسسه صنعتی فولاد، فولادشهر،

aliahoon@cc.iut.ac.ir

## چکیده

تحلیل داده اهمیت خود را در کشف دانش و پیش‌بینی از تصمیم‌گیری در حوزه‌های کاربردی و اطلاعاتی متفاوت به اثبات رسانیده است. تحلیل داده‌های عظیم چالشی جدی در خصوص منابع سخت‌افزاری و نرم‌افزاری مورد نیاز را مطرح نموده است. امروزه تکنولوژی ابر راه‌حلی امید بخش برای این چالش‌ها به وسیله امکان تأمین منابع محاسباتی به صورت مقیاس‌پذیر و فراگیر پیشنهاد کرده است. با این حال چالش‌های بسیاری باقی مانده است، مانند در دسترس بودن نرم‌افزار تحلیلی مورد نیاز برای حوزه‌های کاربردی گوناگون، تخمین و به اشتراک‌گذاری منابع مورد نیاز برای کار تحلیل یا جریان کاری، مدیریت داده در ابر، و طراحی، ارزیابی و اجرای جریان‌های کاری تحلیل، که باید به اینها پرداخته شود. در این مقاله روشی برای بهبود معماری مفهومی تحلیل مبتنی بر ابر به عنوان CLAAaS<sup>۱</sup> که به وسیله آن روند تحلیل با توجه به این که یک معماری همه‌جانبه برای تمامی زمینه‌ها می‌باشد، تسریع بخشیده می‌شود، معرفی شده است.

## کلمات کلیدی

تحلیل داده‌های عظیم، ابر، CLAAaS، هادوپ<sup>۲</sup>

<sup>۱</sup> Cloud-based Analytics-as-a-Service

<sup>۲</sup> Hadoop

تئوری و روش‌های مدیریت سیستم‌های پشتیبان تصمیم‌گیری مبتنی بر سرویس (SODSS) پشتیبانی می‌کند، ارائه داده‌اند[۱]. بخشی از این مقاله نگاهی به چگونگی تحویل محیط DSS مبتنی بر محصول دارد، و فرصت‌ها و چالش‌های مهندسی DSS مبتنی بر سرویس در ابر را تشریح می‌کند. زمانی که داده، اطلاعات و تحلیل به عنوان سرویس تعریف می‌شوند، دیده می‌شود که مکانیزم‌های اندازه‌گیری قدیمی که اساساً زمان‌بر و هزینه‌بر می‌باشند، به خوبی کار نمی‌کنند. سازمان‌ها نیاز دارند که ارزش سطح سرویس و کیفیت را به علاوه هزینه و مدت زمان سرویس تحویل داده شده، در نظر بگیرند. DSS در ابر مقیاس، چشم‌انداز و سرعت اقتصاد را توانمند می‌سازد. این مقاله متشکل از دانش جدید در علم سرویس به وسیله اتصال جنبه‌های استراتژی فناوری اطلاعات به جنبه‌های علم طراحی و پایگاه داده برای مخاطبان وسیع‌تر، می‌باشد[۱]. زالکرنا و همکارانش نیز برای مقابله با چالش داده‌های عظیم، یک طبقه‌بندی برای سیستم‌های جریان کاری تحلیل برای مشخص کردن ویژگی‌های مهم در سیستم‌های موجود ارائه داده‌اند. در این مقاله بر اساس طبقه‌بندی و مطالعه بر روی سیستم‌ها و نرم‌افزارهای تحلیلی موجود، یک معماری مفهومی از تحلیل مبتنی بر ابر به عنوان سرویس (CLAaaS)، بستری برای ایجاد سرویس تحلیل داده‌های عظیم در ابر پیشنهاد شده است. در این تحقیق ویژگی‌های مهم برای CLAaaS به عنوان یک سیستم تأمین سرویس مثل کمک و سفارشی‌سازی برای حوزه و کاربرد خاص، همکاری، معماری ماژولار برای توسعه مقیاس‌پذیر و توافق سطح سرویس (SLA) مطرح شده است[۲].

در کنار چالش مطرح شده به وسیله رشد سریع میزان داده‌ها، این فرصت بسیار بزرگی برای جهانی است که در حال دیجیتالی شدن بیشتر و بیشتر می‌باشد، که اجازه تجمع و تحلیل داده در زمینه‌ای خاص را می‌دهد[۱]. یکی از زمینه‌هایی که تحلیل داده‌های عظیم در آن مطرح می‌باشد و با چالش مربوط به آن مواجه است، علم هواشناسی می‌باشد. اشنیز<sup>۳</sup> و همکارانش، در مقاله خودشان به چالش‌های داده‌های عظیم در علم هواشناسی پرداخته‌اند، حرکت این مقاله به سمت مفهوم تحلیل آب و هوا به عنوان یک سرویس (CAaaS) می‌باشد. این مقاله به دلیل دانش کسب شده در تعامل با داده‌های عظیمی که در نهایت منجر به تولید فواید اجتماعی می‌شود، بر روی تجزیه و تحلیل تمرکز کرده است. اشنیز و همکارانش بر روی CAaaS تمرکز کرده‌اند زیرا بر این باورند که این موضوع راه‌های مفیدی از تفکر درباره مسئله را فراهم می‌کند از جمله: تخصصی کردن مفهوم فرآیند کسب و کار به عنوان یک سرویس که توسعه‌ای توانمند شده از IaaS، PaaS و SaaS به وسیله محاسبات ابری می‌باشد. محاسبات ابری نقش مهمی را در این چارچوب دارد، با این حال فقط به عنوان یک عنصر در مجموعه‌ای از قابلیت‌ها دیده می‌شود که برای ارائه تحلیل آب و هوا به عنوان یک سرویس ضروری می‌باشد. این عناصر به این دلیل ضروری‌اند که در مجموع آن‌ها منجر به تولید ظرفیتی برای خود تجمعی<sup>۴</sup> می‌شوند، که محققان این مقاله احساس می‌کنند که کلید حل بسیاری از چالش‌های داده‌های عظیم در این حوزه می‌باشد. سرویس تجزیه و تحلیل MERRA (MERRA/AS) نمونه‌ای از یک ابر CAaaS می‌باشد که بر اساس این اصل ساخته شده است[۸].

امروزه در جهان کسب‌وکار بسیار پیچیده، سازمان‌ها باید راه‌های مبتکرانه برای متفاوت بودن با رقبا را به وسیله متحد شدن بیشتر، صحت بیشتر، وفق دادن، هماهنگی و چابک‌تر شدن پیدا کنند. آن‌ها نیاز دارند که قادر به پاسخ-گویی به نیازهای کسب‌وکار و تغییر سریع باشند. آن‌ها باید متوجه باشند که صاحب چه داده‌هایی هستند و چگونه متفاوت از دیگران از آن‌ها استفاده کنند. داده‌ها و اطلاعات در حال تبدیل شدن به دارایی اصلی برای بسیاری از سازمان‌ها می‌باشند[1]. تجزیه و تحلیل داده اثر بالقوه خود را در فراهم کردن پشتیبان تصمیم‌گیری در بخش‌های مالی، مدیریتی و علمی با ایجاد قابلیت محاسبات پیچیده برای بدست آوردن دانش، بینش و نتایج تجربی برای اکتشافات علمی اثبات کرده است[2]. واژه تحلیل حوزه وسیعی را که شامل تکنیک‌های تحلیل و پشتیبان تصمیم‌گیری است، ارائه می‌کند. برخی از کارهای تحلیلی قطعی‌اند و دارای خروجی مشخصی می‌باشند که این خروجی‌ها به سیستم پشتیبان تصمیم‌گیری داده می‌شوند ولی برخی دیگر بیشتر حالت اکتشافی دارند به طوری که نتیجه نهایی ممکن است مفید باشد یا نباشد و معمولاً نیاز است که چندین بار مورد بازبینی قرار بگیرند تا نتیجه قطعی حاصل شود. برنامه‌های تحلیل اکتشافی اصولاً پویا هستند. یک فرآیند تحلیل داده که جریان کاری تحلیل نامیده می‌شود[3,4]، به طور کلی شامل وظایفی چون جمع و پاکسازی داده و یا کارهای اکتشافی مثل تعریف و اجرای مدل‌های یادگیری ماشین یا پرسش‌های تحلیلی ساده می‌باشد[۲].

چالش‌هایی نیز در رابطه با داده‌ها برای سازمان‌ها وجود دارد. برای مثال چالش مدیریت میزان زیاد داده‌ها (داده‌های عظیم) که به خاطر انبارهای کوچک‌تر و تکامل تجهیزات جمع‌آوری اطلاعات و داده‌های دیجیتال مثل گوشی‌های همراه، لب‌تاپ‌ها و سنسورها به طور فزاینده در حال بزرگتر شدن می‌باشند[۱]. در نتیجه میزان داده‌های مورد نیاز برای تحلیل با نرخی نمایی در حال رشد می‌باشد و متخصصان تکنولوژی تحلیل داده مانند داده‌کاوی و یادگیری ماشین دانش کافی در مورد این که کدام یک از داده‌ها نیاز است که تحلیل شوند، ندارند[۲]. در سال‌های اخیر داده‌های عظیم به موضوعی فراگیر در کسب‌وکار، علوم کامپیوتر، مطالعات اطلاعات، آمار و بسیاری از زمینه‌ها گشته است[۵]. بر اساس تعریف گروه کارتر، داده‌های عظیم به طور کلی به حجم بالا، سرعت بالای تولید داده و تنوع بالای سرمایه‌های اطلاعاتی اطلاق می‌شود که برای پردازش مبتکرانه و مقرون به صرفه به منظور بهبود ادراک و تصمیم‌گیری مورد استفاده قرار می‌گیرند[۶]. تفکر مبتنی بر سرویس یکی از الگوهای در حال رشد در فناوری اطلاعات می‌باشد که با بسیاری از سیستم‌های دیگر نظیر حسابداری، مالی و عملیاتی در ارتباط می‌باشد[۷]. امروزه تکنولوژی ابر راه‌حلی امید بخش برای این چالش‌ها به وسیله امکان تأمین منابع محاسباتی به صورت مقیاس‌پذیر و فراگیر پیشنهاد کرده است. با این حال چالش‌های بسیاری باقی مانده است، مانند در دسترس بودن نرم‌افزار تحلیلی مورد نیاز برای حوزه‌های کاربردی گوناگون، تخمین و به اشتراک-گذاری منابع مورد نیاز برای کار تحلیل یا جریان کاری، مدیریت داده در ابر، و طراحی، ارزیابی و اجرای جریان‌های کاری تحلیل، که باید به اینها پرداخته شود[۲].

در این زمینه دلن و همکارش، سیستم‌های پشتیبان تصمیم‌گیری مبتنی بر سرویس (DSS در ابر) را مطرح کرده و یک مدل مفهومی که از ارزیابی

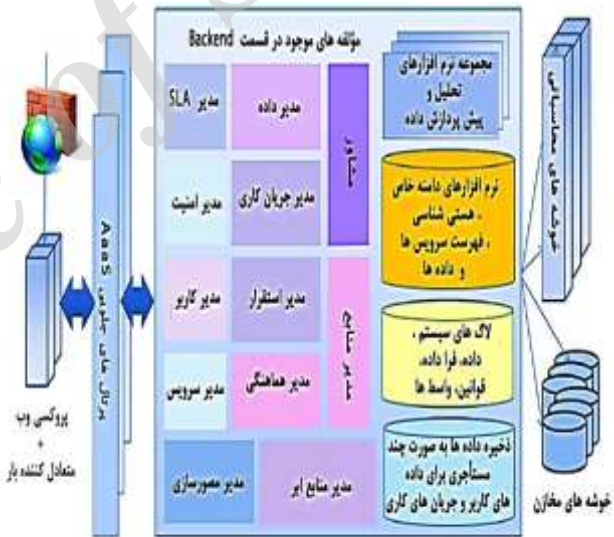
<sup>۳</sup> Schnase

<sup>۴</sup> Self-assembly

سازماندهی مقاله به این صورت است که در بخش دوم معماری مفهومی تحلیل مبتنی بر ابر به عنوان سرویس (CLAAaaS) شرح داده شده است. در بخش سوم تکنولوژی هادوپ و کاهش نگاهت معرفی شده و در بخش چهارم روشی برای بهبود استفاده از هادوپ در CLAAaaS ارائه گشته است و در نهایت، در بخش پنجم، نتیجه‌ی مختصری از موارد بیان شده ارائه می‌گردد.

## ۲- معماری مفهومی تحلیل مبتنی بر ابر به عنوان سرویس (CLAAaaS)

در معماری مفهومی CLAAaaS که در شکل (۱) آمده است، همانند هر سرویس ابر دیگر ترافیک اینترنت از دیواره آتش گذشته و به وسیله متعادل کننده بار و پروکسی وب هدایت و توزیع می‌شود. CLAAaaS دارای دو قسمت واسط کاربر در جلو<sup>۵</sup> که شامل مجموعه‌ای از پرتال‌های وب سفارشی سازی شده برای دامنه‌های متفاوت و گروه‌های کاربری می‌باشد و قسمت موجود در پشت واسط کاربر<sup>۶</sup> که پرتال‌های محاسبات هسته‌ای در آن قرار دارند، می‌باشد. در قسمت راست نیز خوشه‌های محاسباتی و خوشه‌های مخازن که مخازن مقیاس‌پذیری می‌باشند که در هر زمانی لازم باشد، مورد استفاده قرار می‌گیرند، وجود دارند. مؤلفه‌های CLAAaaS در سه بخش مدیریت سرویس، مدیریت جریان کاری و مدیریت داده دسته‌بندی می‌شوند.



شکل (۱): مدل مفهومی CLAAaaS [۲]

(۱) مؤلفه‌های مدیریت سرویس: مدیر توافق سطح سرویس (SLA)<sup>۷</sup>، مدیر امنیت، مدیر کاربر، و مدیر سرویس که وظایف مدیریتی سرویس را انجام می‌دهند. مدیر سرویس مسئول کارکرد و وضعیت کلی سیستم می‌باشد. نظارت بر درخواست‌ها، نگهداری اطلاعات نشست هماهنگی کارکردهای مؤلفه‌های مختلف به وسیله ارتباط کافی با آن‌ها نیز از وظایف این بخش می‌باشد. مدیر SLA مذاکرت SLA را برای سطوح مختلف سرویس‌ها و گروه‌های کاربری را پشتیبانی می‌کند و بر ایجاد SLAها نظارت دارد. مدیر کاربر به ثبت نام

کاربران، نگهداری حساب‌های کاربران و تعریف سطح دسترسی گروه‌های کاربری مختلف کمک می‌نماید. امنیت سیستم توسط مدیر امنیت بر اساس سیاست‌های امنیتی از پیش تعریف شده و سیاست‌های خاص گروه‌های کاربری تأمین می‌شود.

(۲) مؤلفه‌های مدیریت جریان کاری: مدیر داده، مدیر جریان کاری، مدیر مصورسازی<sup>۸</sup>، مدیر هماهنگ‌سازی، مشاور، مدیر منبع و مدیر منابع ابر که همگی با هم تشکیل سیستم مدیریت جریان کار را می‌دهند.

(۳) مؤلفه‌های مدیریت داده: مجموعه نرم‌افزارهای تجزیه و تحلیل داده و سه منبع داده که در کل همگی وظایف پردازش و مدیریت داده را بر عهده دارند. مجموعه نرم‌افزاری می‌تواند شامل پردازش اختصاصی و سفارشی داده‌ها یا نرم‌افزار کاربردی پرس و جو، شامل داده‌های خارجی و یا سرویس‌های نرم‌افزاری باشد. سه منبع داده اصلی در پلتفرم شامل موارد زیر می‌باشند:

- داده‌های مخصوص دامنه: شامل ابزار تحلیلی یک دامنه خاص، اتصال به سرویس‌ها، قالب‌های جریان کاری، هستی‌شناسی، و منابع داده که فهرستی از تسهیلات استفاده مجدد و پرس و جو در آن وجود دارد.
- داده‌های مخصوص سیستم: شامل سیستم‌های کلی مرتبط با داده، فرا داده، لیستی از منابع ابر و کاربران جاری، سیاست‌ها و قوانین از پیش تعریف شده، لاگ‌ها عملیات سرویس و پرتال‌ها می‌باشد.
- داده‌های مخصوص کاربر: شامل داده‌های شخصی کاربران مثل کاربران ثبت نام شده یا SLAهای هر شخص، داده‌های بازگذاری شده یا لینک‌های داده تحلیلی، مدل‌های مصورسازی، مدل‌های تحلیلی و جریان‌های کاری می‌باشد.

## ۳- هادوپ<sup>۹</sup> و کاهش نگاهت<sup>۱۰</sup>

هادوپ یک ابتکار NoSql است، که قابلیت پشتیبانی از این حجم عظیم داده را دارا می‌باشد. ایجاد هادوپ به سال ۲۰۰۵ باز می‌گردد، که به عنوان ابتکاری از یاهو می‌باشد که از تکنولوژی کاهش نگاهت که قبلاً توسط گوگل به کار گرفته شده، استفاده می‌کند. از سال ۲۰۰۹ اکثر شرکت‌ها همچون گوگل، یاهو و فیس بوک برای جستجو، شاخص‌ها و سازماندهی از هادوپ استفاده می‌کنند. هادوپ دارای فایل سیستمی توزیع شده به نام HDFS<sup>۱۱</sup> که در آن پکیج‌های نرم‌افزاری و ابزار گوناگونی وجود دارند که که داده‌های گوناگون از اکثر منابع از لاگ‌ها گرفته تا پایگاه داده‌های رابطه‌ای در آن ذخیره می‌شوند<sup>۱۲</sup>. کاهش نگاهت روشی است که از محاسبات توزیع شده و پردازش موازی داده‌های عظیم را به وسیله کامپیوترهای بسیار قوی پشتیبانی می‌کنند. کاهش نگاهت محاسبات توزیع شده مجموعه بزرگی از داده‌ها را با تعداد زیادی از کامپیوترها (گره‌ها) انجام می‌دهد. در عملیات نگاهت گره اصلی<sup>۱۲</sup> (ریشه) ورودی را گرفته و به زیر مسئله‌های کوچکتر تقسیم کرده و هر زیر گره نیز تا جای ممکن این کار را تکرار می‌کند و در نهایت جواب‌ها به کاهش دهنده پاس داده شده و در این قسمت جواب‌ها با هم ترکیب شده و خروجی

<sup>۸</sup> Visualization

<sup>۹</sup> hadoop

<sup>۱۰</sup> MapReduce

<sup>۱۱</sup> Hadoop Distributed File System

<sup>۱۲</sup> head

<sup>۵</sup> Frontend

<sup>۶</sup> Backend

<sup>۷</sup> Service Level Agreement

## ۵- نتیجه گیری

تجزیه و تحلیل برای ایجاد یک دید بر روی داده‌های عظیم برای تصمیم‌گیری مؤثر تعیین کننده می‌باشد. در این مقاله به معرفی معماری CLAAaaS پرداخته شد زیرا این معماری از جریان‌های کاری تحلیلی برای ساختن پردازش‌های داده گوناگون و مرحله تحلیلی مورد نیاز برای بدست آوردن ارزش از داده‌های استخراج شده می‌پردازد. سپس ساختار قسمت تحلیل این معماری را با مجزا کردن داده‌های متنی عظیم، داده‌های غیر متنی عظیم و داده‌های معمولی که نیاز به استفاده از ظرفیت‌های ویژه ایجاد شده ندارند، بهبود داده شد.

## ۶- منابع

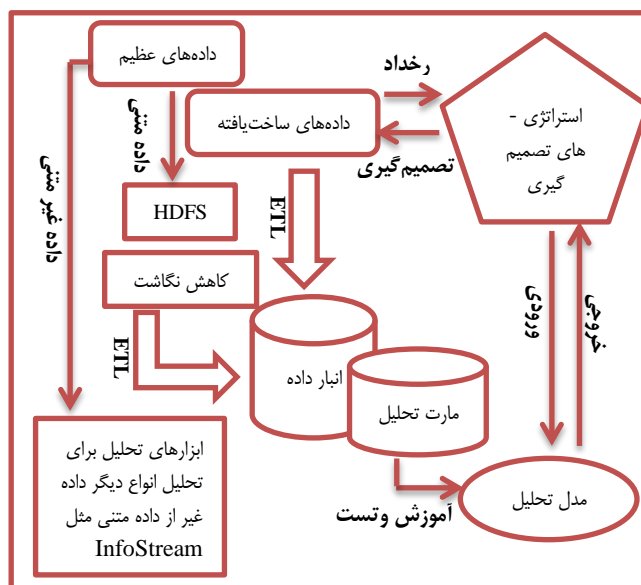
- [1] Haluk Demirkan , Dursun Delen” *Leveraging the capabilities of service-oriented decision support systems: Putting analytics and big data in cloud*”, Decision Support Systems, Volume 55, Issue 1, April 2013, Pages 412-421
- [2] Zulkernine. F, Martin. P, Ying Zou, Bauer. M “*Towards Cloud-based Analytics-as-a-Service (CLAAaaS) for Big Data Analytics in the Cloud*”, 2013 IEEE International Congress on Big Data
- [3] Kim, H., Cho, I., and Yeom, H., 2008. “*A Task Pipelining Framework for e-Science Workflow Management Systems*” In IEEE International Symposium on Cluster Computing and the Grid, (CCGRID'08) pp. 657-662, IEEE.
- [4] myExperiment workflow hosting site: available March 5, 2013 at: <http://www.myexperiment.org/home>.
- [5] Dylan Maltby “*Big Data Analytics*”, ASIST 2011, October 9-13, 2011, New Orleans, LA, USA.
- [6] Nauman Sheikh “*Big Data, Hadoop, and Cloud Computing*”, Implementing Analytics, 2013, Pages 185-197
- [7] H. Demirkan, R.J. Kauffman, J.A. Vayghan, H.-G. Fill, D. Karagiannis, P.P. Maglio, “*Service-oriented technology and management: perspectives on research and practice for the coming decade*”, The Electronic Commerce Research and Applications Journal 7 (4) (2008) 356-376.
- [8] John L. Schnase, Daniel Q. Duffy , Glenn S. Tamkin , Denis Nadeau , John H. Thompson, Cristina M. Grieg, Mark A. McInerney, William P. Webster “*MERRA Analytic Services: Meeting the Big Data challenges of climate science through cloud-enabled Climate Analytics-as-a-Service*” Computers, Environment and Urban Systems, In Press, Corrected Proof, Available online 31 January 2014

به دست می‌آید. یکی از مشکلات موجود در مورد هادوپ این است، تکنولوژی کاهش نگاشت موجود در آن فقط در مورد داده‌های متنی کاربرد دارد و بنابراین در استفاده از آن برای برخی از کاربردهای علوم که دارای داده‌های پیچیده چند بعدی و دودویی هستند، دارای محدودیت می‌باشیم [۸].

## ۴- روش پیشنهادی برای بهبود معماری CLAAaaS

در مقاله [۶] از هادوپ به عنوان موتور برای ETL استفاده شده است. ایده این روش تجمیع و یا ساخت متغیرهای عملکردی از داده‌های عظیم با استفاده از هادوپ می‌باشد در حالیکه داده‌های سنتی که نیاز به متعادل شدن و ترکیب شدن دارند، یک مسیر ETL معمولی را می‌پیمایند. یک مرتبه که داده‌های بزرگ به اندازه‌های قابل مدیریت کاهش می‌یابند، می‌توان با آن‌ها به عنوان داده‌های ساخت یافته معمولی که می‌توانند ذخیره شوند و در پایگاه داده‌ها پردازش شوند برخورد کرد. در این سناریو هادوپ شبیه به یک تکیه‌گاه و یک مؤلفه ETL کارا عمل می‌کند.

با توجه به اینکه در معماری CLAAaaS از انواع نرم افزارها همچون هادوپ و اینفوستریم<sup>۱۳</sup> برای تحلیل داده‌های عظیم استفاده می‌شود و همچنین به دلیل اینکه در برخی از زمینه‌ها ممکن است داده‌های هنوز به میزانی نرسیده باشند که بتوان به آن‌ها داده‌های عظیم گفت، مدل موجود در شکل (۲) برای ساخت یافته‌تر کردن قسمت تحلیلی معماری CLAAaaS پیشنهاد شده است. در این پیشنهاد مدل معرفی شده توسط نعمان شیخ [۶] با اضافه کردن شرط-های داده‌های متنی و داده‌های غیر متنی به قسمت تحلیلی معماری CLAAaaS اضافه می‌شود که با توجه به اینکه هادوپ یک ابزار قدرتمند NoSql می‌باشد در تمامی زمینه‌های داده‌های متنی به کار گرفته می‌شود و همچنین ظرفیت‌های موجود در مورد داده‌های معمولی که نیاز به اینگونه پردازش‌ها ندارند به کار گرفته نمی‌شوند و این موضوع موجب تسریع و بهبود روند در پردازش موازی و توزیع شده داده‌های عظیم می‌شود.



شکل (۲): روش ارائه شده برای بهبود معماری CLAAaaS

<sup>۱۳</sup> InfoStream