

## مروری بر تکنیک‌ها و ابزارهای پردازش جریان در کلان داده

<sup>1</sup> رویا جعفری، <sup>2</sup> اصغر درویشی

<sup>1</sup> دانشجوی دکتری تخصصی نرم افزار دانشگاه آزاد اسلامی واحد تهران جنوب، Jafarir.ir@gmail.com

<sup>2</sup> استادیار و عضو هیئت علمی گروه مهندسی نرم افزار دانشگاه آزاد اسلامی واحد تهران جنوب، A\_Darvishi@azad.ac.ir

### چکیده

با رشد روزافزون دستگاه‌های دیجیتال و ارتباطات شبکه‌ای گسترده میان آنها بطور ویژه در دستگاه‌های اینترنت اشیا و همچنین شبکه‌های اجتماعی در دهه‌های اخیر، حجم بزرگی از داده در حال تولید می‌باشد که کلان داده نامیده می‌شود. کلان داده را می‌توان توالی نامحدود و پیوسته‌ای از داده در نظر گرفت که با سرعت در حال تولید هستند و به نام جریان داده می‌باشند. جریانی که اگر در زمان مناسب پردازش شود، منافع زیادی برای بسیاری در پی خواهد داشت و در غیراینصورت، داده مورد نظر ارزش خود را از دست می‌دهد. درکنار نگرانی‌هایی مانند چگونگی ذخیره این حجم از داده و نیز در دسترس بودن آن، پردازش به موقع آن، چالشی مهم برای بسیاری از سازمان‌ها می‌باشد. روش‌های پردازش داده به دو صورت دسته‌ای و مبتنی بر جریان می‌باشند که راه حلی برای پردازش حجم زیادی داده ایستا که از قبیل در سیستم ذخیره شده اند می‌باشد بدین صورت که داده پس از جمع‌آوری از منابع مختلف تولیدکننده داده و ذخیره و نگهداری در سیستم، بصورت دسته‌ای مورد پردازش قرار می‌گیرد که انجام آن زمانبر است ولی در روش‌های مبتنی بر جریان، داده به محض دریافت و بصورت بلادرنگ پردازش می‌شود و همین امر موجب افزایش سرعت پردازش و در دسترس بودن داده را فراهم می‌نماید. پردازش جریان داده، نیاز به روش‌ها و ابزارهایی برای دریافت داده و تجزیه و تحلیل بر آن دارد که راهکارهایی برای حل چالش‌های پردازش جریان پیوسته داده می‌باشد. در مطالعه پیش‌رو، روش‌ها و تکنیک‌های کلان داده و فریمورک‌های محبوب که سرعت پردازش داده را افزایش می‌دهند مورد بررسی قرار گرفته است. رویکرد پردازش جریان برای زبان طبیعی، روش‌های توزیع شده پردازش جریان داده، موتورهای پردازش و در نهایت، یک دسته‌بندی سیستماتیک صورت گرفته درباره چالش‌ها و راهکارهای آن‌ها در حوزه پردازش جریان داده برای انتخاب فریمورک مناسب‌تر ارائه شده است.

### واژه‌های کلیدی:

کلان داده، پردازش جریان، پردازش دسته‌ای، بلادرنگ، موتورهای پردازش

### ۱- مقدمه

در دهه‌های اخیر بویژه دهه جاری، جهان مورد تهاجم همه انواع تکنولوژی در بخش‌های مختلف می‌باشد [4]. دیجیتالی شدن هر حوزه‌ای از زندگی بشر منجر به انفجاری در تولید داده شده است [5] که شامل افزایش پیوسته حجم و جزئیات داده‌های جذب شده توسط سازمان‌ها از قبیل ظهور رسانه‌های اجتماعی، اینترنت اشیا و مالتی مدیا که منجر به ایجاد یک جریان بسیار عظیم داده در شکل‌های مختلفی مانند ساختاریافته، نیمه ساختاریافته و ساختاریافته شده است [6] خلق داده با سرعت بی سابقه‌ای در حال ایجاد است که آن را کلان داده می‌نامیم. کلان داده‌ها به هر مجموعه بسیار وسیع و پیچیده‌ای اطلاق می‌شود که پردازش آن با استفاده از برنامه‌های پردازش سنتی مشکل خواهد بود. کلان داده‌ها مجموعه‌ای از تکنیک‌ها و فناوری‌هایی است که نیازمند فرم‌های جدیدی برای آشکار نمودن مقادیر پنهان بر روی مجموعه داده‌های متنوع، پیچیده با مقیاس بزرگ می‌باشند [1]. در واقع قابلیت استخراج اطلاعات مفید از مجموعه داده‌های بزرگ و یا جریان‌های داده بدلیل حجم، تنوع و سرعت داده‌ها است، امکانی که قبلاً وجود نداشت [۱۳] که یکی از تکنیک‌های رایج استخراج داده، خوشه بندی می‌باشد، مجموعه‌ای از خوشه‌های محدود شده که اغلب روی تعداد خوشه‌ها تعریف می‌شود [۱۴].

کلان داده را می‌توان متشکل از توالی نامحدود و بزرگی از عناصر داده‌ای که به صورت پیوسته و با سرعت زیاد در حال تولید می‌باشند نامید که به آن جریان داده می‌گوییم. بدلیل توالی نامحدود جریان کلان داده و سرعت تولید آن، امکان ذخیره سازی تمام داده‌های جریان وجود ندارد. داده‌ها با

سرعت زیاد و به صورت بلادرنگ در حال ایجاد هستند. در بعضی کاربردها نیاز آن است که به محض وارد شدن داده به درخواست کاربر پاسخ داده شود و ممکن است نتوانیم به اندازه کافی صبر کنیم تا یک گزارش در سیستم برای مدت طولانی پردازش شود [2].

تا سال ۲۰۱۹، ۲۶ میلیارد از دستگاه‌های اینترنت اشیا بهم متصل بودند و انتظار می‌رود که این رقم تا ۲۰۲۵ به ۸۰ میلیارد دستگاه در جهان برسد. سیستم‌های پردازش جریان به سرعت در حال تبدیل شدن به عنصر اصلی پشته پردازش داده در بسیاری از شرکت‌ها می‌شوند [۱۵]. این در حالی است که الگوریتم‌های سنتی پردازش جریان بدلیل استفاده بسیار زیاد از منابع (مانند حافظه) برای پردازش پویا و توزیع شده داده قادر به انجام فعالیت در مقابل این جریان داده وسیع نیستند. عناوینی مانند: تکامل جریان داده، زمان اجرا، مصرف حافظه، چند بعدی بودن وسیع، مفاهیم جهش، تاخیر و عدم تعادل از جمله چالش‌های الگوریتم‌های سنتی در مواجهه با نوع جدید جریان داده می‌باشند [4]. همانگونه که تکنولوژی‌های پردازش جریان در مسیر بلوغ قدم می‌گذارند، سازمان‌های بیشتری جذب این بخش شده و بر روی تحول دیجیتال سرمایه‌گذاری می‌کنند. مجموعه داده‌های جهانی و بدون مرز مانند وب لاگ‌ها، آمار استفاده از تلفن‌های همراه و شبکه‌های حسگر روز به روز در حال افزایش است و برنامه‌های جدید تجزیه و تحلیل جریان در طیف گسترده‌ای از صنایع، شناسایی و پیاده‌سازی می‌شوند. جریان داده بسیاری از پلت فرم‌های شبکه‌های اجتماعی مانند: فیس‌بوک، توئیتر و لینکدین شاهرگ زندگی آنها می‌باشد و لذا استخراج معنی دار و آگاهی به موقع از داده‌های نامحدود، چالش پیش رو در بیشتر بخش‌ها می‌باشد [7]. پردازش جریان داده به عنوان کلیدی در حوزه پردازش جریان اطلاعات است و هدف آن درک رفتار داده و استخراج اطلاعات برای تصمیم‌گیری بهتر می‌باشد [8]. در ابتدا پردازش داده‌ها بصورت دسته‌ای بود که بعد از جمع‌آوری داده‌ها از منابع گوناگون و انجام عملیات استخراج، انتقال و بارگذاری<sup>۱</sup> در مخزنی از داده، با توجه به نوع کسب و کار و نیازهای سازمان شروع به ساخت و ایجاد گزارش‌های داشبوردی برای سطوح مختلف مدیران صورت می‌پذیرفت که به این مخزن انباره داده می‌گویند که شامل پردازش داده‌های ایستا می‌باشد. در حال حاضر، جریان داده بصورت پویا و پیوسته است که نیاز به پردازش لحظه‌ای را ایجاد می‌کند. پردازش‌هایی با کمترین تاخیر، معنی دار و مقیاس‌پذیر که کمبودهای سیستم‌های پردازش دسته‌ای را جبران کنند و نیاز کاربران را در لحظه، پاسخگو باشند، به اصطلاح بلادرنگ و در زمان باشند. از این دست پردازش‌های بلادرنگ می‌توان به سناریوهایی مانند کشف تقلب در تراکنش‌های مالی و یا تجزیه و تحلیل سامانه‌های مراقبت سلامت و بهداشتی که شامل حسگرهای دیجیتال و اینترنت اشیا هستند نام برد. به دلیل آنکه انواع خاصی از جریان‌های داده مانند ارزش سهام، تراکنش‌های کارت اعتباری، ترافیک داده‌های حساس به زمان بیمار و روند اخبار در صورتی که فوری پردازش نشوند، ارزش خود را از دست خواهند داد لذا توانایی مدیریت و پردازش جریان‌های مداوم و پیوسته داده‌ها به یکی از اصلی‌ترین بخش‌ها در سازمان‌های داده محور امروزی تبدیل شده است [7].

آنچه پیش روی این بررسی می‌باشد بدین صورت خواهد بود: در بخش ۲ مروری بر تکنیک‌های کلان داده خواهیم داشت. رویکرد جریان برای پردازش زبان طبیعی را در بخش ۳ بیان می‌کنیم. در بخش ۴، مروری بر پردازشگرهای بلادرنگ خواهیم داشت. فریمورک‌های توزیع شده پردازش جریان داده در بخش ۵ بیان شده‌اند و دسته‌بندی سیستماتیک برای انتخاب بهتر ابزارهای پردازش جریان داده را در بخش ۶ ارائه نموده ایم و در نهایت، در بخش ۷ به بیان خلاصه‌ای از آنچه در این مطالعه مورد بحث و بررسی قرار گرفت، می‌پردازیم.

## ۲- مروری بر تکنیک‌های کلان داده

در [6] تمرکز بر مرور و بررسی تکنیک‌هایی است که برای پردازش و به کارگیری، نمایش اطلاعات و تحلیل آنها بکار گرفته شده است. این تکنیک‌ها مجموعه‌ای از قواعد را در برمیگیرد که شامل داده‌های آماری، یادگیری ماشین، شبکه‌های عصبی، تحلیل شبکه‌های اجتماعی، پردازش سیگنال، تشخیص الگو، روش‌های بهینه‌سازی و روش‌های مجازی‌سازی می‌باشد که در ادامه به آنها اشاره می‌شود:

روش‌های بهینه‌سازی: روش‌هایی که برای برطرف نمودن و حل مسائل تا حدودی در زمینه‌هایی مانند فیزیک، زیست، مهندسی و اقتصاد بکار گرفته شده است. شبیه‌سازی تیرید، الگوریتم ژنتیک و حرارت دهی کوانتومی.

روش‌های آماری: تکنیک‌های آماری برای بدست آوردن همبستگی‌ها و روابط نسبی بین اهداف گوناگون بکار می‌رود. الگوریتم تقریبی برای رگرسیون یکنواخت چند متغیره در مقیاس بالا و الگوریتم‌های آماری موزی.

داده کاوی: تحلیل خوشه بندی، دسته بندی، رگرسیون و یادگیری قانون.

یادگیری ماشین: کشف دانش و تصمیم‌گیری هوشمندانه خودکار از شاخصه‌های این بخش می‌باشد. SVM موزی و شبکه‌های عصبی مصنوعی.

روش تجسم: این روش تکنیک‌های مورد استفاده برای جداول، تصاویر، نمودارها و در کل، نمایش بصری برای درک داده می‌باشد.

تحلیل شبکه‌های اجتماعی: یک تکنیک کلیدی است که در جامعه‌شناسی مدرن ظهور پیدا کرده است و به روابط اجتماعی بر حسب نظریه تئوری شبکه می‌نگرد که شامل گره‌ها و روابط بین آن‌ها است.

<sup>1</sup> ETL = Extract, Transfer, Load

### ۳- رویکرد جریان برای پردازش زبان طبیعی

در [10]، یک معماری جدید توزیعی و تکنولوژی جدید برای مقیاس پذیری تحلیل متن با حفظ سرعت نرخ رشد جریان های داده ارائه شده است و زنجیره ای کامل از پردازشگرهای متنی بر روی چندین ماشین مجازی می باشد. همچنین ماژول های پردازش زبان طبیعی که شامل ماشین های مجازی می باشند و بیشتر آنها از ابزارهای خط لوله IXA وارد می شوند، پیشنهاد شده است. جدول ۱ ماژول های نصب شده در ماشین مجازی زبان انگلیسی را نشان می دهد. داخل ماشین های مجازی با استفاده از فریمورک Storm برای هر محاسبه جریان مدیریت می شوند و هر ماژول پردازش زبان از گره های bolt داخل توپولوژی Storm تشکیل شده است. هنگامی که تاپل جدیدی وارد می شود گره های bolt یک دستور خارجی را برای ارسال محتوای تاپل به یک جریان ورودی استاندارد فرا می خواند. خروجی ماژول پردازش جریان از خروجی جریان استاندارد دریافت شده به گره بعدی در توپولوژی ارسال می شود. بنابراین هر ماژوی یک سند NAF<sup>T</sup> با سند و افزودن خط لوله برای هر ابزار جدید برای نوشتن و خواندن NAF را دریافت می کند. در نتیجه برای دانش ما جامع ترین خط لوله NLP موجود است. این روش می تواند برای تحلیل متون فارسی مورد بهره برداری قرار گیرد. از دو مجموعه داده اطلاعات خبری راجع به Car و Wikinews برای آزمایش ماژول ها استفاده شده است.

جدول ۱. ماژول های پردازش زبان نصب شده بر روی ماشین های مجازی

ستون آخر نشان دهنده نسخه خط لوله ای است که ماژول ها استفاده می کنند

ماژول	توضیحات	خط لوله
<i>ixa-pipe-tak</i>	Tokenizer, sentences splitter	(۲و۱)
<i>ixa-pipe-pos</i>	POS tagger	(۲و۱)
<i>ixa-pipe-parse</i>	Constituency parser	(۲)
<i>TimePro</i>	Time expression recognition	(۲و۱)
<i>ixa-pipe-nerc</i>	Named Entity Recognition	(۲و۱)
<i>WSD</i>	Word Sense Disambiguation	(۲و۱)
<i>dbpedia-spotlight</i>	Named Entity Disambiguation	(۲و۱)
<i>ixa-pipe-coref</i>	Coreference resolution	(۲)
<i>MATE</i>	Dependency parser and Semantic Role Labeling	(۲و۱)
<i>opinion miner</i>	Opinion detection and Opinion holders to targets	(۲)
<i>Factuality</i>	Factuality	(۲و۱)
<i>eCoref</i>	Event coreference	(۲و۱)

در مجموعه داده Car ماژول MATE-SRL تقریباً ۸۰ درصد زمان پردازش را به خود اختصاص می دهد. در مجموعه داده Wikinews ماژول MATE-SRL ۵۶ درصد از کل زمان را می گیرد و در ادامه، *ixa-pipe-coref* ۲۱ درصد از زمان را مصرف می کند. این نتایج پیشنهاد می دهند که هر دو ماژول کاندیدهای مناسبی برای موازی سازی هستند. برای موازی سازی ماژول ها از یک خط لوله NLP که شامل ۴ ماژول: *ixa-pipe-tok*، *ixa-pipe-pos*، *ixa-pipe-nerc* و UKB و ابزاری برای ابهام زدایی از حس کلمه بر اساس گراف است، استفاده شده است.

جدول ۳. عملکرد خط لوله پردازش زبان طبیعی با تنظیمات متفاوت

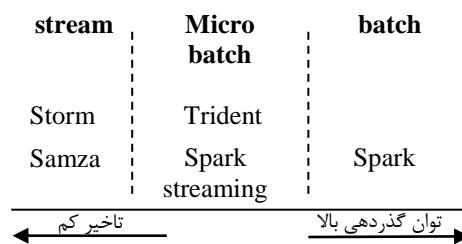
Storm دارای یک نمونه و Storm<sub>۲</sub> دارای دو نمونه و Storm<sub>۴</sub> دارای چهار نمونه از ماژول WSD است

	زمان کل	کلمات(ثانیه)	فرستاده شده(ثانیه)	بهره (/)
<i>۱۰۰ documents</i>				
Pipeline	۲۱m ۱۶s	۱۰۸,۸	۴,۲	-
Storm	۱۸m ۴۳s	۱۲۳,۵	۴,۸	۱۲,۰
Storm <sub>۲</sub>	۱۰m ۴۸s	۲۱۴,۳	۸,۴	۴۹,۳
Storm <sub>۴</sub>	۷m ۴۶s	۲۹۷,۶	۱۱,۶	۶۳,۵
Storm <sub>۸</sub>	۷m ۴۴s	۲۹۹,۱	۱۱,۷	۶۳,۷
Storm <sub>۱۶</sub>	۷m ۴۸s	۲۹۶,۱	۱۱,۶	۶۳,۳
<i>۱۰۰۰ documents</i>				
Pipeline	۳h ۱۵m ۱۶s	۱۰۱,۲	۴,۲	-
Storm	۲h ۵۰m ۲۱s	۱۱۶,۰	۴,۸	۱۲,۸
Storm <sub>۲</sub>	۱h ۴۰m ۳۷s	۱۹۶,۵	۸,۱	۴۸,۵
Storm <sub>۴</sub>	۱h ۱۴m ۲۵s	۲۶۵,۶	۱۰,۹	۶۱,۹
Storm <sub>۸</sub>	۱h ۱۰m ۴۵s	۲۷۹,۳	۱۱,۵	۶۳,۸
Storm <sub>۱۶</sub>	۱h ۱۱m ۳۷s	۲۷۶,۰	۱۱,۳	۶۳,۳

جدول ۳ نشان دهنده اجرای خط لوله NLP در تنظیمات مختلف است: ۱۰۰ سند و ۱۰۰۰ سند. با اجرای چندین نمونه از UKB عملکرد کلی افزایش معناداری پیدا می کند. بالاترین بهره برای WSD در ۵ نمونه است که برابر 63.7 درصد است و با افزایش نمونه ها، بهره افزایش نمی کند. نتایج نشان دهنده فضای بزرگی برای بهبود پردازش NLP است.

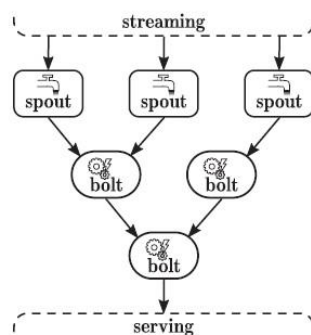
۴- پردازشگرهای بلادرنگ

در [11] اشاره می‌شود که پردازشگرهای جریان با وجود مشترک بودن مفاهیم و اصل کار خود تفاوت اساسی بین سیستم‌هایی که به سرعت قابل دستیابی پردازش تبدیل می‌شوند مانند تاخیر وجود دارد که در شکل ۲ نشان داده شده است. پردازش سریع داده‌ها به محض ورود، منجر به کاهش تاخیر خواهد شد در حالیکه بافر کردن و پردازش دسته‌ای داده‌ها با وجود آنکه کارایی را افزایش می‌دهد اما زمان سپری شده داده در خط لوله را افزایش می‌دهد. سیستم‌های جریان محور مانند Storm و Samza تاخیر بسیار کم و هزینه نسبتاً بالا را ایجاد می‌کنند در حالیکه سیستم‌های دسته‌ای به منابع کارایی غیرموازی دست پیدا می‌کنند که هزینه تاخیر آنها بصورت غیرقابل اجتنابی برای برنامه‌های بلادرنگ، بالا می‌باشد. فاصله میان تاخیر کم و گذردهی بالا بسیار زیاد است و بعضی از سیستم‌ها مانند Storm Trident و Spark Streaming استراتژی‌های دسته‌های کوچک<sup>۳</sup> دارند که برای معاوضه مابین تاخیر و گذردهی بکار می‌روند: در Storm Trident، تاپل‌ها به هدف ماندن در یک مدل پردازش زمانی به سود افزایش گذردهی، به دسته‌هایی تبدیل می‌شوند، در حالیکه در Spark Streaming اندازه دسته‌های داده با هدف ایجاد تاخیر کم، محدود می‌باشد.



شکل ۲. انتخاب یک مدل پردازش، موازنه‌ای بین گذردهی و تاخیر

Storm: اولین سیستم پردازش جریان توزیع شده که معروف به هادوپ بلادرنگ<sup>۴</sup> می‌باشد. یک خط لوله داده در Storm، توپولوژی نامیده می‌شود که همانطور که در شکل ۳ مشخص می‌باشد گراف جهت دار جریان داده می‌باشد و لبه‌های جهت دار میان گره‌ها مراحل پردازش را نشان می‌دهند: گره‌های spout، گره‌هایی هستند که داده از آنجا وارد می‌شود و تاپل‌های ورودی را به گره‌های جریان‌های پایینی تر که bolt نامیده می‌شوند و عملیات پردازش روی تاپل‌های ورودی را انجام می‌دهند، منتشر می‌کنند. Storm ضمانتی برای ترتیب پردازش تاپل‌ها نمی‌دهد و همچنین حداقل یکبار پردازش را از طریق ردیابی تاپل‌های وارد شده به توپولوژی فراهم می‌کند. یک مکانیسم فشار برگشتی<sup>۵</sup> برای داده وارد شده معرفی نمود که داده‌هایی که سرعت ورود آنها بیشتر از آن است که Storm بتواند آنها را پردازش کند، داده مجدد وارد گره‌های توپولوژی به هدف پردازش گردد. بدون چنین مکانیسمی، پردازش می‌تواند تبدیل به گلوگاهی شود که موجب کاهش توان گذردهی خواهد شد.



شکل ۳. جریان داده در توپولوژی Storm: داده از لایه جریان وارد و پس از عبور از مولفه‌های Storm

در خروجی به لایه خدمت می‌رسد

Storm Trident، توپولوژی‌ها، گراف‌های جهت دار بدون چرخه می‌باشند و از چرخه پشتیبانی نمی‌کنند که مناسب پیاده‌سازی الگوریتم‌های تکرارپذیر نمی‌باشد. Trident با تاپل‌های مجزا کار نمی‌کند بلکه بصورت دسته‌های کوچک با معرفی اندازه دسته‌ها به عنوان پارامتری برای افزایش گذردهی و کاهش تاخیر می‌باشد (تاخیرهایی در حد چند میلی‌ثانیه). همه دسته‌ها با ترتیب و بصورت

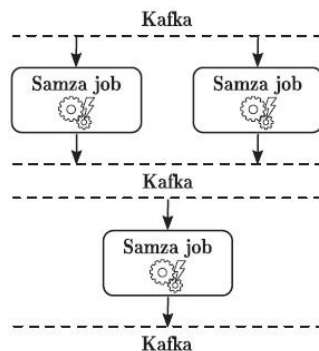
<sup>3</sup> Micro-batching

<sup>4</sup> Hadoop of real-time

<sup>5</sup> Backpressure mechanism

<sup>6</sup> DAG = Directed Acyclic Graph

متوالی، یکی بعد از دیگری، پردازش می‌شوند. همینطور امکان پردازش چندین دسته بصورت موازی نیز وجود دارد. در بالای Trident, Storm برنامه‌ای برای مدیریت خط‌پذیری با دقت یکبار پردازش معنا فراهم می‌کند. Samza: یک پردازشگر جریان با مدل پردازش در زمان و حداقل یکبار پردازش معنا می‌باشد. جریان‌ها تقسیم می‌شوند و پیام‌ها (مانند آیتم‌های داده) در هر قسمتی، مرتب می‌شوند در حالیکه هیچ ترتیبی بین پیام‌ها در قسمت‌های مختلف وجود ندارد. برخلاف Storm, Samza فقط به ZooKeeper وابسته نیست بلکه در بالای Hadoop YARN برای خط‌پذیری اجرا می‌شود.



شکل ۴. جریان داده در یک خط لوله تجزیه و تحلیل معمولی Samza

همانطور که در شکل ۴ مشخص می‌باشد، یک کار Samza مراحل پردازش را در خط لوله تجزیه و تحلیل نشان می‌دهد که بطور تقریبی با یک گره bolt در توپولوژی Storm مطابقت دارد. خروجی‌های تولید شده در هر کار Samza همیشه در کافکا نوشته می‌شود زیرا ممکن است مورد نیاز یک کار Samza دیگر باشد که ورودی خود را از کافکا دریافت می‌کند. این طراحی مراحل پردازش را جدا می‌کند و در نتیجه توسعه به سادگی صورت می‌پذیرد. از دیگر مزایای این پردازشگر جریان، بافر کردن داده بین گام‌های پردازش است که می‌تواند توسط تیم‌های دیگر همان سازمان مورد استفاده قرار گیرد و در نهایت، حذف نیاز به الگوریتم فشار پستی را منجر خواهد شد زیرا مشکلی در پرشدن موقتی یک کار خاص بدلیل اندازه منطقی توسعه با کافکا، وجود ندارد. چون Samza پیام‌ها را به ترتیب پردازش می‌کند و نتایج را بصورت ماندگار نگهداری می‌نماید از گم شدن داده جلوگیری می‌نماید.

Spark Streaming: یک فریمورک پردازش جریان که به عنوان جانشین غیررسمی هادوپ ذکر شده است چون در مقایسه، مزایای متعددی را ارائه می‌دهد که مهمترین آن، بهبود عملکرد قابل توجه از طریق ذخیره سازی در حافظه می‌باشد. به ویژه الگوریتم‌های تکراری (الگوریتم‌های یادگیری ماشین مانند k-means یا رگرسیون لجستیک) به ترتیب اندازه، شتاب گرفته‌اند زیرا لزومی به نوشتن و بارگذاری داده از دیسک در هر مرحله از پردازش نیست. باتوجه به مزایای این عملکرد، Spark الگوریتم‌های متنوع یادگیری ماشین را از طریق کتابخانه‌های MLlib فراهم می‌نماید. مانند Storm و Samza خط‌پذیر عمل می‌کند و از مقیاس پذیری منابع اختصاص داده شده بصورت پویا پشتیبانی می‌کند. برای جلوگیری از دست دادن داده حتی برای منابع داده‌ای غیر قابل اطمینان، Spark Streaming اجازه پیش نویس داده<sup>۷</sup> از جایی که داده می‌تواند مجدد بعد از خطا اجرا شود را می‌دهد. جدول ۴ مقایسه‌ای بین فریمورک‌های اشاره شده را بیان می‌کند.

جدول ۴. مقایسه Apache Spark Streaming و LinkedIn's Samza, Apache Storm/Trident

	Storm	Trident	Samza	Spark Streaming
Strictest guarantee	At-least-once	Exactly-once	At-least-once	Exactly-once
Achievable latency	<< ۱۰۰ ms	< ۱۰۰ms	< ۱۰۰ms	< ۱s
State management	Yes	Yes (small state)	Yes	Yes
Processing model	One-at-a-time	Micro-batch	One-at-a-time	Micro-batch
Backpressure mechanism	Yes	Yes	Not required(buffering)	Yes
Ordering guarantees	No	Between batches	Within stream partitions	Between batches
Elasticity	Yes	Yes	No	Yes

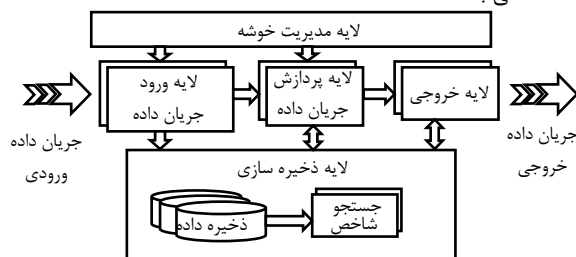
<sup>7</sup> WAL = Write-Ahead Log

## ۵- فریمورک‌های توزیع شده پردازش جریان داده

در [7]، ۲ محور اصلی برای پردازش جریان ذکر شده است: سیستم‌های مدرن پردازش جریان<sup>۸</sup> و موتورهای پردازش جریان داده<sup>۹</sup> که به توضیح آن در ادامه خواهیم پرداخت.

### ۵-۱- سیستم‌های پردازش جریان داده

این سیستم‌ها سعی در ترکیب قابلیت‌های پردازش دسته‌ای و جریان در یک یا چند خط لوله پردازش موازی داده را دارند. رویکرد سنتی پردازش داده، پردازش دسته‌ای می‌باشد، مثلا داده‌های تراکنش‌های مالی که در یک انباره داده تجمیع و ذخیره می‌گردد و سپس، داده‌های ذخیره شده مورد پردازش و تحلیل قرار می‌گیرند و یا سیستم جستجوی وب که با جستجوی مداوم در وب به دنبال استخراج دانش می‌باشد و پردازش آن بصورت دسته‌ای است تا شاخص‌های جستجو را در هر ساعت بروزرسانی نماید. مشکل هر دو سناریو، تأخیر زیاد بین جمع‌آوری داده و استفاده از آن‌ها را نشان می‌دهد. هدف این سیستم‌ها غلبه بر تأخیر با پردازش بلافاصله داده‌ها در زمان رسیدن با در نظر گرفتن مقیاس‌پذیری و خط‌پذیری است. ورودی DSPS‌ها، جریان‌های پیوسته داده هستند و هر DSPS علاوه بر تفاوت در منطق کسب و کار و دامنه تحت پوشش، دارای یک خط لوله پردازش داده عمومی می‌باشد که از ۵ لایه شامل: لایه دریافت داده، لایه پردازش جریان داده، لایه ذخیره‌سازی، لایه مدیریت منبع و لایه خروجی تشکیل شده است که در شکل ۵ قابل مشاهده می‌باشد.



شکل ۵. معماری سیستم پردازش جریان داده

لایه دریافت داده: استفاده از سرویس‌های پیام‌رسان از مدل‌های سنتی مانند MQTT، RabbitMQ و ActiveMQ تا محصولات جدیدتر مانند ZeroMQ و NSQ می‌باشد.

لایه پردازش جریان داده: محلی که برنامه‌های پردازش جریان یا کارها اجرا می‌شوند. موتورهای پیشرفته مدیریت پردازش<sup>۱۰</sup> شامل Samza، KSQL و SQL stream Blaze می‌باشند.

لایه ذخیره‌سازی: راه‌حل‌های ذخیره‌سازی داده برای پشتیبانی از معماری DSPS از سیستم‌های فایل سنتی مانند HDFS و BFS<sup>۱۱</sup> تا دیتابیس‌های رابطه‌ای فایل توزیع شده مانند PostgreSQL، دیتابیس‌های درون‌حافظه‌ای مانند VoltDB، مبتنی بر سند مانند MongoDB، سیستم‌های ذخیره‌سازی مبتنی بر گراف مانند Neo4j، دیتابیس‌های NoSQL مانند Cassandra و NewSQL مانند CockroachDB می‌باشند.

لایه مدیریت منبع: فرایندهای موجود در یک سیستم توزیع شده می‌تواند از طریق شبکه یا برخی از فضاهای ذخیره‌سازی مشترک خواندن/نوشتن، پیام را مستقیماً تبادل نماید. مدیریت اعضای گروه و مدیریت فراداده هنگام ساختن یک خوشه توزیع شده از بخش‌های مهم این لایه می‌باشد. مدیریت منابع پیشرفته مانند ZooKeeper، Mesos، Tupperware و سرویس مدیریت شده‌ای مانند Borg می‌باشند.

لایه خروجی: نتایج حاصل شده از خط لوله پردازش جریان می‌تواند به برنامه‌های دیگری هدایت شوند یا وارد ابزار تجسم‌سازی گردند و یا یک داشبورد نظارتی/هشدار را ایجاد نمایند (شکل ۵).

ابزارهای تجسم‌سازی به چهار دسته زیر تقسیم می‌شوند:

- تصویرسازی گراف با استفاده از ابزاری مانند TreeNetViz.
- تصویرسازی متنی که شامل تصویرسازی متن ایستا با ابزاری مانند WordClouds یا PhraseNet و یا متن پویا با ابزاری مانند SparkClouds یا TextFlow.
- تصویرسازی نقشه برای داده‌های جغرافیایی با ابزاری مانند BirdVis.
- تصویرسازی چندمتغیره داده برای انواع داده‌های عمومی با استفاده از ابزارهای مقیاس‌پذیری چند بعدی.

<sup>8</sup> DSPTS = Data Stream Processing System

<sup>9</sup> DSPE = Data Stream Processing Engine

<sup>10</sup> DSME = Data Stream Management Engine

<sup>11</sup> Baido File System

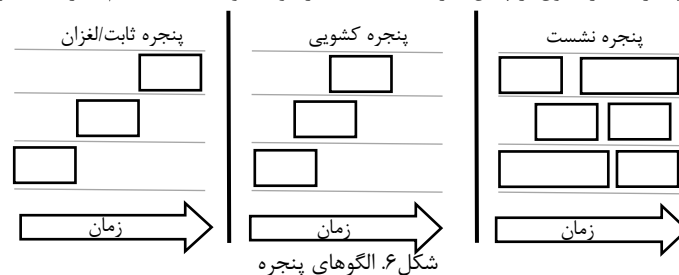
۵-۲- موتور پردازش جریان داده

موتورهای پردازش جریان شامل خطوط پردازش جریان داده را به عنوان یک DAG<sup>۱۲</sup> منطقی متصل شده به کارهای پردازش جریان می باشد که DAG را اجرا می کند بدین صورت که آیتم های جریان را از یک منبع دریافت می کند، روی آن پردازش انجام می دهد و نتایج را بر می گرداند و از چندین مولفه تشکیل شده است که شامل منابع داده و سینک، درایور برنامه، مدیر جریان و پردازشگرهای جریان است. موتورهای پردازش جریان توزیع شده با تمرکز بر مقیاس پذیری حجم کلان داده بصورت پیوسته طراحی شده اند [۱۸].

جدول ۵. مقایسه موتورهای پردازش جریان داده

معیار	Storm	Spark Streaming	Flink	Kafka Streams	IBM Streams
Version	2.0	2.4.3	1.8	2.3.0	4.3.0
Model	Native	Micro batch, continuous processing	Native	Native	Native
Partitioning	Shuffle, field, partial key, all, global, none, direct, local, custom	Hash, range	Same, hash, random	Same, hash, round-robin	Same, hash, random, round-robin, entire, modulus, range, DB2, auto
State management	Key value, redis	Write-ahead and other state stores	Key-value store	Key-value, hash map	Tuple history, operator windows, operator custom logic, primitive operators
Message delivery	Atleast-once, Exactly-once through Trident	Exactly- and At-least-once	Exactly-once	Exactly-once	Exactly-once
Fault tolerance	Checkpoint, Stream replay	Checkpoint	Stream replay, checkpoint	Stream replay	Checkpoint, consistent region, replay

داشتن یک دید کلی از جریان داده های ورودی، امکان پذیر نیست لذا داده ها بصورت قطعه هایی با پنجره های زمانی محدود در بالای یک جریان داده بی انتها، پردازش می شوند. پنجره های مبتنی بر زمان به سه گروه تقسیم می شوند(شکل ۶): پنجره های ثابت: جریان داده ورودی را به سگمنت هایی با اندازه ثابت که هر کدام یک طول پنجره، زمان شروع و زمان پایان دارند، مشخص می شود. پنجره های کشویی: جریان داده ورودی را به سگمنت هایی با اندازه ثابت تقسیم می کند که هر کدام یک طول پنجره، زمان شروع و زمان پایان دارند. اگر بازه کشویی به اندازه طول پنجره باشد، تبدیل به پنجره ثابت می شود در نتیجه یک یا چند بخش از داده شامل بیش از یک پنجره کشویی می گردد. پنجره های نشست: بر خلاف دو پنجره دیگر، طول از پیش تعریف شده ای ندارد و معمولاً با فاصله عدم تحرک که بزرگتر از حد آستانه است تعریف می شود.



این موتورها به چهار دسته زیر تقسیم می شوند:

- ۱- سیستم های بر اساس پرس و جو مانند AsterixDB و NiagaraCQ
- ۲- سیستم های یادگیری ماشین توزیع شده آنلاین مانند SAMOA<sup>۱۳</sup>
- ۳- سیستم های تجزیه و تحلیل گراف جریان مانند GraphJat
- ۴- فریمورک های همه منظوره پردازش جریان داده مانند Flink و Spark Streaming

۶- چالش ها و راهکارها برای پردازش بلادرنگ جریان کلان داده

در [6]. به بررسی سیستماتیک ادبیات پردازش جریان که به گردآوری مطالعات مرتبط با پردازش جریان بلادرنگ و چالش ها و تحولات انجام شده است، می پردازد. از میان ۶۷۹،۵۴۷ مقاله انتخاب شده برای استخراج داده، ۷۴ مقاله بعد از ارزیابی کیفی برای تحقیق باقی ماندند که در این بررسی

<sup>12</sup> Directed Acyclic Graph

<sup>13</sup> Scalable Advanced Massive Online Analysis

به پیاده سازی چالش‌ها، همراه با رویکردهای توسعه‌ای انباره داده‌های بلادرنگ و سیستم‌های پردازش جریان کلان داده و مقایسه میان آن‌ها را فراهم نموده است. در متدولوژی تحقیق شامل سه فاز: برنامه، هدایت و بررسی گزارش می‌باشد. شکل ۷ متدولوژی تحقیق و فرایندهای جستجو برای فعالیت‌های تحقیق را نشان می‌دهد.



شکل ۷. متدولوژی تحقیق

یک فرایند کاملا ساختاریافته که شامل موارد زیر می‌باشند، دنبال شده است: اهداف تحقیق، مشخص نمودن سوالات تحقیق، سازماندهی جستجوهای پایگاه داده، مطالعه انتخاب‌ها، غربالگری مطالعات مرتبط، استخراج داده، سنتز نتایج و نهایی نمودن گزارشات مورد بررسی.

#### ۶-۱- اهداف تحقیق

شامل : - توسعه دادن عناوین مرتبط کتابخانه به پردازش جریان بلادرنگ در طول فاز ETL یا سایر فازها و ایجاد یک مجموعه داده.  
- مشخص نمودن کارهای مرتبط که به چالش‌های پردازش جریان بلادرنگ، ETL و انباره داده بلادرنگ جهت می‌دهد.  
- تشخیص شکاف‌های تحقیق.  
- مشخص کردن رویکردها و راهکارهای موجود برای این چالش‌ها در زمان پیاده سازی پردازش جریان بلادرنگ برای داده‌های ناهمگون، ساختاریافته/غیرساختاریافته و روشن نمودن شباهت‌ها و تفاوت‌های میان آنها.  
آنچه اهمیت دارد آن است که بتوان سوالات تحقیق را فرموله نمود، پرسش‌هایی که در جدول ۲ قرار گرفته است.

جدول ۶. پرسش‌های تحقیق

پرسشهای تحقیق	عبارات پرسش‌های تحقیق	انگیزه
RQ 1	کدامیک انتشارات مرتبط برای تحقیق پردازش جریان بلادرنگ می‌باشند؟	مشخص نمودن آنکه تحقیق پردازش جریان بلادرنگ می‌تواند هدف مناسبی برای مطالعات آینده باشد با ارزیابی کیفیت کانال‌های هدف بر اساس رتبه بندی آن‌ها و دسته بندی انتشارات برای دوره زمانی انتخاب شده، نوع و مناطق جغرافیایی تحقیقات.
RQ 1.1	چه تعداد عنوان در سال‌های بین ۲۰۱۱ تا ۲۰۲۰ منتشر شده است؟	
RQ 1.2	کدام کانال‌ها و مناطق جغرافیایی تحقیقات پردازش جریان بلادرنگ را هدف قرار می‌دهند؟	
RQ 1.3	چه نوع تحقیقاتی، رویکردهایی و برنامه‌هایی برای مطالعه انتخاب می‌شوند؟	
RQ 1.4	چه تعداد از رویکردهای مطالعات بصورت تجربی تایید شده است؟	
RQ 2	با کدام چالش‌ها در طول پیاده سازی پردازش جریان داده‌ها مواجه شده‌اند؟	با ارزیابی و ساده سازی پردازش جریان بلادرنگ با نیازمندی‌های متفاوت و به منظور تشویق نمودن برای جنبه‌های مختلف با هدف تحقیقات در آینده. شناسایی برنامه‌هایی که پردازش جریان بلادرنگ انجام می‌دهند و به چیزی بجز انباره داده نیاز دارند
RQ 2.1	کدام ابزارها/تکنولوژی‌ها/ رویکردها برای حل این چالش‌ها توسعه داده شده است؟	
RQ 2.2	کدام پیاده سازی در حوزه اینترنت اشیا و رسانه‌های اجتماعی با پردازش بلادرنگ جریان داده گزارش شده است؟	
RQ 3	در زمان پردازش جریان بلادرنگ برای انباره داده بلادرنگ، چه رویکردها/ابزارهایی برای حل چالش‌ها در مرحله ETL گزارش می‌شود؟	برای تشخیص رویکردهای متصل گزارش شده داده جریان/نیمه جریان، ساختاریافته/نیمه ساختاریافته موجود در ادبیات پردازش جریان بلادرنگ، چالش‌هایی در پیاده سازی انباره داده بلادرنگ شناسایی می‌شود.
RQ 3.1	کدام شکل و ساختار داده برای پیاده سازی رویکردها استفاده می‌شود؟	
RQ 4	در حل چالش‌های پردازش جریان بلادرنگ چه اتفاقاتی گزارش می‌شوند؟	ارزیابی اتفاقات فراهم شده بوسیله این مطالعات، منجر به رویکردهای دقیقتری می‌شود و نیز تشخیص مجموعه داده‌های استفاده شده در هنگام آزمایش رویکرد پیشنهاد شده (مصنوعی یا واقعی). تشخیص اثر رویکردها بر روی عملکرد و هزینه و ارزیابی فرکانس تطبیق معیار عملکردی خاص برای اعتبارسنجی.
RQ 4.1	چه مجموعه داده‌هایی برای آزمایش بوسیله این رویکردها گزارش شده است؟	
RQ 4.2	چه معیار عملکردی برای آزمایش این رویکردها با مجموعه داده‌های واقعی سازگار می‌باشد؟	



- فرایند هدایت نمودن این بررسی شامل چهار مرحله می باشد که به این ترتیب می باشند:
- ۱- مطالعات اصلی مرتبط از رایج ترین کتابخانه های دیجیتال جستجو شده اند.
  - ۲- انتخاب مطالعات بر اساس معیارهای استخراج/شامل بودن انجام شده است.
  - ۳- سپس، معیار ارزیابی کیفیت برای بدست آوردن کیفیت بیشتر صورت می پذیرد.
  - ۴- جستجوهای مرتبط برای استخراج مقالات کاندید در طول این گام انجام می شود.

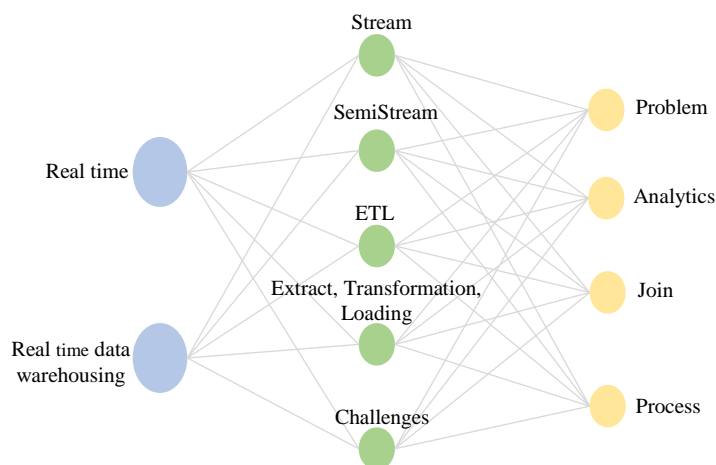
#### ۲-۶- تحقیق نیمه خودکار

تحقیق نیمه خودکار در ۷ کتابخانه دیجیتال زیر انجام پذیرفته که در جدول ۷ آمده است.  
جدول ۷. کتابخانه های دیجیتال مورد استفاده

کتابخانه های دیجیتال
ACM Digital Library [ <a href="http://dl.acm.org">http://dl.acm.org</a> ]
IEEE eXplore [ <a href="http://ieeexplore.ieee.org">http://ieeexplore.ieee.org</a> ]
ScienceDirect [ <a href="https://www.sciencedirect.com">https://www.sciencedirect.com</a> ]
IGI Global [ <a href="https://www.igi-global.com/search/">https://www.igi-global.com/search/</a> ]
Inderscience Online [ <a href="https://www.inderscienceonline.com/">https://www.inderscienceonline.com/</a> ]
Hindawi [ <a href="https://www.hindawi.com/">https://www.hindawi.com/</a> ]
MDPI (Multidisciplinary Digital Publishing Institute) [ <a href="https://www.mdpi.com/">https://www.mdpi.com/</a> ]
arXiv [ <a href="https://arxiv.org/search/cs">https://arxiv.org/search/cs</a> ]
Taylor&Francis Online [ <a href="https://www.tandfonline.com/">https://www.tandfonline.com/</a> ]

علاوه بر موارد بیان شده در جدول ۷، چند کتابخانه دیجیتال دیگر نیز که شامل محدودیت های دسترسی بودند، بررسی شدند. مرتب سازی های ممکن از جمله تحقیق را می توان در شکل ۸ مورد توجه قرار داد. کلیدهای اصلی به عنوان کلید شناسایی برای تحقیق پردازش جریان بلادرنگ انتخاب شدند و کلیدهای اصلی انتخاب شده با کلیدهای ثانویه یا کایدهای تجمیعی، مرتبط هستند. ترکیب این کلیدها، عملگرهای بولین و کارکترهای جایگزین<sup>۱۴</sup>، رشته جستجوی نهایی ایجاد می شود:

(real time OR real tiem data warehous\*) AND  
(stream OR semi-stream OR ETL OR challenges OR (extract\*,transform\*,load\*)) AND  
(join OR process\* OR analy\* OR problem\*)



شکل ۸. جستجوی کلمات کلیدی با استفاده از پایگاه دانش برای شناسایی آثار

جدول ۸، جملات تحقیق نهایی را نشان می دهد که برای جستجو در کتابخانه های دیجیتال بیان شده مورد استفاده قرار گرفته اند. انتخاب بر اساس شامل بودن در محدوده و خارج بودن از محدوده تحقیق می باشد. شامل بودن را مواردی مانند بودن در محدوده پردازش بلادرنگ جریان داده، هدف

پرسش‌های تحقیق، نتایج شامل مقالات ترویجی باشند و می‌بایست در مورد توسعه و برنامه‌های پردازش بلادرنگ بحث و بررسی نموده باشد می‌توان نام برد. همچنین انگلیسی نبودن و مرتبط با مباحث پردازش بلادرنگ جریان در انباره داده یا کلان داده‌ها، خارج از محدوده تحقیق می‌باشند. همچنین پرسش‌های تحقیق مورد ارزیابی قرار گرفته‌اند. برای مثال، در پاسخ به پرسش اول در جدول ۲ که شامل: کدام انتشارات مرتبط با تحقیق پردازش جریان بلادرنگ می‌باشند؟ باید گفت تحلیل وجود توسعه و چالش‌ها برای پردازش بلادرنگ جریان داده، چالش کلیدی برای توسعه تکنولوژی‌های هوش تجاری است. تشخیص و شناسایی انتشارات با کیفیت بالا و تحلیل علم‌سنجی بر اساس فرااطلاعات در ناحیه پردازش بلادرنگ جریان مورد نیاز می‌باشد.

جدول ۸. استراتژی‌های جستجو برای کتابخانه‌های دیجیتال

فیلترهای اعمال شده	جمله جستجو	کتابخانه دیجیتال
[(01/01/2011) TO 01/31/2020]	[[All: real time] OR [All: real time data warehous*]] AND [[All: stream] OR [All: etl] OR [All: challenges] OR [All: extract*,transform*,load*]] AND [[All: join] OR [All: process*] OR [All: analy*] OR [All: problem*]]	ACM digital library
2011- Jan2020	(real time OR real time data warehouse*) AND (stream OR semi-stream OR ETL OR challenges OR (extract*,transform*,load*)) AND (join OR process* OR analytics OR problems)	IEEEExplore
2011- Jan2020	(real time OR real time data warehouse) AND (stream OR semi-stream OR ETL OR challenges OR (extraction,transformation,loading)) AND (join OR process OR analysis OR problems)	ScienceDirect
Computer Science, (2011- Jan2020)	^(real time OR real time data warehouse*) AND (stream OR semi-stream OR ETL OR challenges OR (extract*,transform*,load*)) AND (join OR process* OR analy* OR problem*)^	SpringerLink
Individual Journal Articles (2011- Jan2020)	(real time) AND (data warehousing or stream processing)	IGI Global
2011- Jan2020	(real time OR real time data warehouse*) AND (stream OR semi-stream OR ETL OR challenges OR (extract*,transform*,load*)) AND (join OR process* OR analytics OR problems)	Inderscience Online
2011- Jan2020	"real time stream processing OR real time data warehouse"	Hindawi
2011- Jan2020	"real time stream processing"	MDPI
2011- Jan2020	"real time data warehouse OR real-time stream processing"	arXiv
2011- Jan2020	(real time OR real time data warehouse*) AND (stream OR semi-stream OR ETL OR challenges OR (extract*,transform*,load*)) AND (join OR process* OR analy* OR problem*)	Taylor & Francis Online

## ۷- نتیجه گیری

سیستم‌های دسته‌محور، برای دهه‌ها، وزنه‌های سنگینی بودند اما برای داده‌هایی با ماهیت پیوسته و نامحدود که در بسیاری از کاربردهای دنیای واقعی تولید می‌شوند، کاربردی ندارد. سیستم‌های جریان‌محور یا به عبارتی دیگر، سیستم‌هایی که به محض دریافت داده آنها را پردازش می‌کنند از تناسب طبیعی بیشتری با داده‌های کلانی که همه‌روزه در حال تولید هستند، برخوردار هستند [11]. در این میان نیاز به پلت‌فرم‌هایی برای مدیریت مقدار زیاد داده در زمینه‌های مشخص شده در یک روش کاملاً توزیع شده می‌باشد، رویکردی که داده‌کاو سنتی قادر به انجام آن نبود و به یک گلوگاه عملکردی غیرقابل قبول تبدیل شده بود [۲۰]. در این بررسی به اهمیت پردازش جریان داده و حفظ ارزش‌های پنهانی آن پرداخته شد. مروری بر تکنیک‌های کلان داده داشتیم، در نتیجه تقاضای قابل توجهی برای برنامه‌های کلان وجود دارد [۱۶]. روش‌هایی برای داده‌ها با حجم زیاد و قواعدی که بتوان با داده‌های آماری و یادگیری ماشین و تجزیه و تحلیل الگوها به کشف دانش پرداخت. سپس به نقش پردازش جریان داده در بازاریابی اشاره شد و اینکه تا چه حد ابزارهایی مانند شبکه‌های اجتماعی و داده‌های تولید شده توسط آنها می‌تواند منجر به سودآوری و افزایش وفاداری مشتریان در کسب و کارهای مختلف گردد. نقش جریان داده برای پردازش زبان طبیعی مورد بررسی قرار گرفت. با استفاده از خط لوله پردازش می‌توان به تحلیل متون همچون -فارسی پرداخت. پردازشگرهای بلادرنگ و جریان‌محوری مانند Storm، Samza، Spark Streaming و Strom Trident بررسی شدند و قابلیت‌ها و امکانات و نحوه پردازش جریان داده توسط آنها بیان شد. در ادامه، دو مفهوم اصلی سیستم‌های مدرن پردازش جریان و موتورهای پردازش جریان داده ارائه گردید. مولفه‌های شکل دهنده آنها و نحوه دریافت داده در طول زمان بررسی شدند و در نهایت یک تحقیق ساختارمند که بصورت سیستماتیک به جمع‌آوری مقالات در حوزه پردازش جریان داده‌ها و انبار داده بلادرنگ می‌پردازد، بیان شد. حرکت از اهداف تحقیق آغاز گردید، پرسش‌هایی مرتبط با حوزه مورد نظر برای رسیدن به مقالات و تحقیقات مرتبط ایجاد گردید. با بررسی مقالات موجود در چندین کتابخانه دیجیتال و فیلتر نمودن موارد خارج از محدودی تحقیق به ۷۲ مقاله مرتبط در فاصله بین سال‌های ۲۰۱۱ تا ۲۰۲۰ با

موضوع تحقیق رسیدیم. با رشد نمایی داده‌های جریان و نیاز به پردازش‌های لحظه‌ای و ارسال بازخورد به کاربران، مطمئناً با چالش‌هایی در بخش منابع همچنان مواجه هستیم که این موضوع نیاز به فراهم کردن راه‌های نو و تازه برای استفاده از جریان پیوسته که بصورت لحظه‌ای تولید می‌شود، در آینده‌ای نزدیک خواهد شد.

#### ۸- منابع

۱. کلکلی، منصور و رجایی، امیر، پردازش داده‌های کلان (Big data) ۱۳۹۶، کنفرانس سالانه پارادایم‌های نوین مدیریت در حوزه هوشمندی، تهران.
۲. رضایی، اکرم و حق جو، مصطفی، الگوریتمی برای پردازش سریع جریان کلان داده، ۱۳۹۷، کنفرانس ملی پیشرفت‌های اخیر در مهندسی و علوم نوین، قرچک.
۳. زاهد، عطیه و سخی، محمدرضا، مروری بر تکنیک‌های کلان داده (Big Data)، ۱۳۹۵، سومین کنفرانس ملی مهندسی برق و کامپیوتر سیستم‌های توزیع شده و شبکه‌های هوشمند، کاشان.
4. Bahri, M., Bifet, A., Gama, J., and Gomes, H. M., and Maniu, S. 2021. Data stream analysis: Foundations, major tasks and tools. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 11(3): e1405.
5. Mittal, V., Gangodkar, D., & Pant, B. 2021. Big Data Technologies: A Comprehensive Survey. *Inventive Communication and Computational Technologies*, 103-114.
6. Mehmood, E., & Anees, T. 2020. Challenges and solutions for processing real-time big data stream: A systematic literature review. *IEEE Access*, 8, 119123-119143.
7. Isah, H., Abughofa, T., Mahfuz, S., Ajerla, D., Zulkernine, F., & Khan, S. 2019. A survey of distributed data stream processing frameworks. *IEEE Access*, 7, pp: 154300-154316.
8. Alghushairy, O. S. 2021. A Genetic Algorithm-Based Local Outlier Factor for Efficient Big Data Stream Processing (Doctoral dissertation, University of Idaho).
9. MODI, J. Big Data Stream Processing in Marketing Analytics.
10. Agerri, R., Artola, X., Beloki, Z., Rigau, G., & Soroa, A. 2015. Big data for Natural Language Processing: A streaming approach. *Knowledge-Based Systems*, 79, 36-42.
11. Wingerath, W., Gessert, F., Friedrich, S., & Ritter, N. 2016. Real-time stream processing for Big Data. *IT-Information Technology*, ۱۸(4), 186-194.
12. Dubey, A. K., Gupta, R., & Mishra, S. 2021. Data Stream Clustering for Big Data Sets: A comparative Analysis. In *IOP Conference Series: Materials Science and Engineering* (Vol. 1099, No. 1, 012030).
13. Fan, W., & Bifet, A. 2013. Mining big data: current status, and forecast to the future. *ACM SIGKDD explorations newsletter*, 14(2), 1-5.
14. Abualigah, L., Gandomi, A. H., Elaziz, M. A., Hamad, H. A., Omari, M., Alshinwan, M., & Khasawneh, A. M. 2021. Advances in meta-heuristic optimization algorithms in big data text clustering. *Electronics*, 10(2), 101.
15. Affetti, L., Margara, A., & Cugola, G. 2020. TSpooon: Transactions on a stream processor. *Journal of Parallel and Distributed Computing*, 1۴۰, 65-79.
16. Alwasel, K., Calheiros, R. N., Garg, S., Buyya, R., Pathan, M., Georgakopoulos, D., & Ranjan, R. 2021. BigDataSDNSim: A simulator for analyzing big data applications in software-defined cloud data centers. *Software: Practice and Experience*, ۵۱(5), 893-920.
17. Sahal, R., Breslin, J. G., & Ali, M. I. 2020. Big data and stream processing platforms for Industry 4.0 requirements mapping for a predictive maintenance use case. *Journal of Manufacturing Systems*, ۵۴, 138-151.

18. Henning, S., & Hasselbring, W. 2021. Theodolite: Scalability benchmarking of distributed stream processing engines in microservice architectures. *Big Data Research*, ۲۰, 100209.
19. Pishgoo, B., Azirani, A. A., & Raahemi, B. 2021. A hybrid distributed batch-stream processing approach for anomaly detection. *Information Sciences*, ۵۴۳, 309-327.
20. Esposito, C., Ficco, M., Palmieri, F., & Castiglione, A. 2015. A knowledge-based platform for big data analytics based on publish/subscribe services and stream processing. *Knowledge-Based Systems*, ۱۹, 3-17.