



کلاس‌بندی بسته‌های رمزنگاری شده  
(Encrypted Packet Classification)  
درس یادگیری عمیق

فاطمه مهدوی

استاد درس:

دکتر محمدظاهر پیلهور

بهار ۱۳۹۸

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

## فهرست مطالب

۱.....	چکیده
۱.....	مقدمه
۳.....	کارهای مرتبط
۵.....	مدل پیشنهاد شده
۵.....	مجموعه داده و مشخصاتش [۲]
۶.....	پیش پردازش داده
۷.....	معماری پیشنهادی
۹.....	نتایج
۱۰.....	تحلیل
۱۱.....	مراجع

## چکیده

کلاسه‌بندی ترافیک اینترنت با رشد سریع شبکه فعلی اینترنت و برنامه‌های کاربردی آنلاین اهمیت بیشتری پیدا کرده است. مطالعات متعددی در این زمینه انجام شده، که منجر به ارائه رویکردهای مختلفی شده است. در اکثر این روش‌ها از ویژگی‌های استخراج شده از مجموعه داده توسط یک کارشناس، برای کلاسه‌بندی ترافیک شبکه استفاده می‌کنند. در مقابل، در مرجع [۱]، یک رویکرد مبتنی بر یادگیری عمیق ارائه شده، که در آن هر دو مرحله استخراج ویژگی‌ها و کلاسه‌بندی در یک سیستم ادغام شده است. طرح پیشنهادی آنها، به نام بسته عمیق<sup>۱</sup>، می‌تواند هر دو ویژگی ترافیکی اعم از ویژگی شبکه و شناسایی برنامه‌های کاربردی را استخراج نماید، همچنین می‌تواند ترافیک رمزنگاری شده را شناسایی کند. پس از یک مرحله اولیه پیش پردازش بر روی داده‌ها، بسته‌ها به عنوان ورودی به چارچوب بسته عمیق وارد می‌شوند که از آتوانکدر انباشته و شبکه عصبی کانولوشنال تشکیل شده است. در مرجع [۱] استفاده از شبکه عصبی کانولوشنال نتایج بهتری را نسبت به آتوانکدر انباشته بدست آورده است. در این پروژه به جای شبکه عصبی کانولوشنال، شبکه کپسوله پیشنهاد شده است که با توجه به مزایای آن نسبت به شبکه‌های کانولوشن انتظار می‌رود که نتایج بهتری حاصل گردد.

## مقدمه

کلاسه‌بندی ترافیک شبکه یک وظیفه مهم در شبکه‌های ارتباطی مدرن است. فناوری‌های کلاسه‌بندی ترافیک، به دلیل پیاده‌سازی مکانیزم‌هایی برای افزایش کیفیت خدمات شبکه (QoS)، امنیت، حسابداری، طراحی و مهندسی، نسبت به دهه گذشته، در هر دو بخش دانشگاه و صنعت توجه بسیاری را به خود جلب کرده است. صنعت شبکه و همچنین جامعه محققان تلاش‌های زیادی را برای تحقیق در مورد این فناوری‌ها اختصاص داده و چندین تکنیک کلاسه‌بندی ارائه نموده‌اند. با این حال، گسترش مداوم اینترنت و تکنولوژی‌های تلفن همراه یک محیط پویا ایجاد کرده به طوری که در آن برنامه‌های کاربردی و خدماتی جدید هر روز پدیدار شده و موارد موجود به طور مداوم در حال تحول هستند. همچنین رمزنگاری در اینترنت امروزه در حال گسترش است و به عنوان پایه‌ای برای ارتباطات ایمن کاربرد دارد. بنابراین ایجاد، تکامل و امنیت برنامه‌های کاربردی، باعث می‌شود که کلاسه‌بندی ترافیک، چالش بزرگی در زمینه تحقیقات اینترنتی باشد.

برای نمونه اهمیت کلاسه‌بندی ترافیک شبکه را می‌توان از معماری نامتقارن لینک‌های دسترسی به شبکه‌های جدید، که بر اساس این فرض که کاربران بیش از آنچه که آپلود کنند، دانلود می‌کنند، طراحی شده‌اند، دریافت. با این حال، فراگیر شدن برنامه‌های کاربردی با تقاضای متقارن (از قبیل برنامه‌های کاربردی (P2P)، Voice over IP (VoIP) و تماس ویدیویی)، خواسته‌های کاربران را نسبت به پیش‌فرض ذکر شده، تغییر داده است. بنابراین،

---

<sup>۱</sup> Deep Packet

به منظور ارائه یک تجربه رضایت‌بخش برای کاربر، به دانشی از سطح برنامه‌های کاربردی در تخصیص منابع کافی برای چنین برنامه‌هایی لازم است.

کلاسه‌بندی ترافیک را می‌توان براساس هدف نهایی آن بدین صورت طبقه‌بندی کرد: مرتبط‌سازی ترافیک با رمزنگاری (به عنوان مثال ترافیک رمزنگاری شده)، کپسوله‌سازی پروتکل (به عنوان مثال از طریق VPN یا HTTPS)؛ براساس برنامه‌های کاربردی خاص (به عنوان مثال اسکایپ) یا طبق نوع برنامه (به عنوان مثال جریان و چت)، که همگی جزء مشخصات ترافیک می‌باشند. برخی از برنامه‌های کاربردی (مانند اسکایپ و فیسبوک) از سرویس‌های متعددی مانند چت، تماس صوتی و انتقال فایل پشتیبانی می‌کنند. این برنامه‌های کاربردی هم نیاز به شناسایی برنامه کاربردی و هم کار خاص مرتبط به آن را دارند. تکنیک‌های کلاسه‌بندی ترافیکی بسیار کمی به این روندهای چالشی اشاره نموده‌اند.

حال در مورد برخی از مهم‌ترین چالش‌های کلاسه‌بندی ترافیک شبکه بحث می‌کنیم. در ابتدا، تقاضای روزافزون برای حریم خصوصی و رمزنگاری داده‌ها، میزان ترافیک رمزنگاری شده در اینترنت امروزه را افزایش داده است. پروسه رمزنگاری، داده‌های اصلی را به شکل شبه تصادفی تبدیل می‌کند تا رمزگشایی آن را سخت کند. در عوض، این موضوع باعث می‌شود که داده‌های رمزنگاری شده به ندرت شامل هرگونه الگوهای قابل تشخیص برای شناسایی ترافیک شبکه باشد. در اوایل دهه ۹۰، یک روش ساده و سریع از تکنیک‌های کلاسه‌بندی اولیه ترافیک این بود که، پورت‌های لایه انتقال را با برنامه‌های کاربردی خاص مرتبط می‌نمود. اما دقت پایین و عدم اطمینان آن باعث توسعه روش‌های بازرسی بسته‌های عمیق<sup>۱</sup> (DPI) گردید. رویکرد DPI بسته‌ها را تجزیه و تحلیل کرده و آنها را طبق برخی نشانه یا الگوی ذخیره شده کلاسه‌بندی می‌کند. با این حال، تکنیک‌های DPI نیازمند بررسی مجدد بوده و از نظر محاسباتی مخصوصاً در مورد شبکه با پهنای باند بالا، کارآمد نیستند. علاوه بر این، آنها اغلب به وسیله ترافیک کپسوله شده، رمزنگاری شده یا محو شده که مانع تجزیه و تحلیل بار می‌شوند، منحرف می‌گردند. بنابراین، کلاسه‌بندی دقیق ترافیک رمزنگاری شده در شبکه‌های مدرن به یک چالش تبدیل شده است.

همچنین لازم به ذکر است که بسیاری از کلاسه‌بندی‌های پیشنهادی شبکه ترافیک، مانند بازرسی بارگیری و همچنین یادگیری ماشین و روش‌های آماری، نیاز به الگوها یا ویژگی‌های استخراج شده توسط متخصصین دارند، که این روند دارای خطا، وقت‌گیر و پرهزینه است.

علاوه بر این انتخاب ویژگی‌های موثر و قابل اطمینان برای تجزیه و تحلیل ترافیک هنوز یک چالش جدی است. به طور کلی کلاسه‌بندی ترافیک شبکه به طور عمده به دو دسته تقسیم می‌شود: کلاسه‌بندی مبتنی بر جریان، با استفاده از خواصی مانند بایت جریان بر ثانیه، مدت زمان بر جریان و کلاسه‌بندی مبتنی بر بسته، با استفاده از خواصی مانند اندازه، مدت زمان بین ورود بسته‌ای از بسته‌ی اول (یا nام).

با وجود اینکه مطالعات فراوانی در مورد کلاسه‌بندی ترافیک شبکه وجود دارد، اکثر آنها بر کلاسه‌بندی خانواده پروتکل، با عنوان تعیین مشخصات ترافیک (مانند جریان، چت و P2P)، به جای شناسایی یک برنامه کاربردی

---

<sup>۱</sup> Deep Packet Inspection

واحد که به عنوان شناسایی برنامه شناخته می‌شود (مانند Spotify, Hangouts و Bittorrent) تمرکز دارند. در مقابل، در مرجع [۱] روش بسته عمیق، هم برای تشخیص و هم شناسایی ترافیک شبکه پیشنهاد شده است. مزایای این روش پیشنهادی که آن را نسبت به دیگر طرح‌های کلاسه‌بندی برتر می‌کند، به شرح زیر است:

در روش بسته عمیق، نیازی به متخصص برای استخراج ویژگی‌های مربوط به ترافیک شبکه نیست. با توجه به این رویکرد، مرحله دشوار مربوط به پیدا کردن و استخراج ویژگی‌های متمایز حذف شده است.

روش بسته عمیق می‌تواند ترافیک را در هر دو سطح کلاسه‌بندی (شناسایی برنامه و تعیین مشخصات ترافیک) با نتایج پیشرفته‌ای نسبت به سایر کارهایی که بر روی مجموعه داده‌ی مشابه انجام شده است، شناسایی کند. [۲]

روش بسته عمیق می‌تواند به دقت کلاسه‌بندی یکی از سخت‌ترین کلاس‌های برنامه‌های کاربردی شناخته شده با عنوان P2P را انجام دهد.

در این پروژه مانند مرجع [۱]، ما بر روی تجزیه و تحلیل ترافیک رمزنگاری شده معمولی و ترافیک رمزنگاری شده از طریق یک شبکه خصوصی مجازی (VPN) تمرکز کرده‌ایم. تشخیص ترافیک VPN یک کار چالش برانگیز است که هنوز به طور کامل حل نشده است. تونل‌های VPN برای حفظ حریم خصوصی داده‌های به اشتراک گذاشته شده، با حفظ اتصال فیزیکی شبکه و رمزنگاری در سطح بسته استفاده می‌شود، بنابراین شناسایی برنامه‌های کاربردی در حال اجرا از طریق این سرویس‌های VPN بسیار مشکل است.

در ادامه، ابتدا به بررسی کارهای مرتبط انجام شده در این زمینه پرداخته می‌شود. سپس در زمینه مدل پیشنهاد شده به بررسی مواردی از قبیل نحوه‌ی تولید مجموعه داده [۲] مورد استفاده، پیش پردازش داده و معماری پیشنهادی می‌پردازیم. در نهایت نیز با ارائه نتایج و تحلیل خود، کارهایی که در ادامه این پروژه باید صورت گیرد، را بیان می‌نماییم.

## کارهای مرتبط

مطالعات بر روی اندازه بسته و کلاسه‌بندی ترافیک بر مبنای جریان در اوایل دهه ۹۰ توسط پکسون و همکاران [۴] آغاز شد، جایی که برخی از ویژگی‌های آماری مانند طول بسته، زمان‌های درهم آمیختگی و مدت زمان جریان، برای ردیابی پروتکل‌ها مطلوب فرض شده است. سپس در مراجع متفاوتی از توابع آماری به شکل‌های مختلفی برای به دست آوردن کارایی مطلوب استفاده شده است.

تکنیک‌های تعیین مشخصات ترافیک در تحقیقات فعلی به طور گسترده‌ای مورد توجه نبوده‌اند. علاوه بر این، اکثر آنها بر روی نوع خاصی برنامه کاربردی یا دستگاه تمرکز دارند. وانگ و همکاران [۵] یک مدل برای تعیین مشخصات ترافیک P2P ارائه دادند. آنها ویژگی‌هایی از جریان‌های چندگانه و جریان‌های جمع شده را در خوشه‌ها استخراج کردند تا رفتار برنامه کاربردی از نوع P2P را استخراج کنند. شری و همکاران [۶] یک سیستم DPI را پیشنهاد کردند که می‌تواند بدون نیاز به رمزگشایی بارگیری رمزنگاری شده را بررسی کند، بنابراین حریم خصوصی ارتباطات حفظ می‌شود، اما تنها می‌تواند ترافیک HTTPS را پردازش کند.

تعدادی از روش‌های کلاسه‌بندی یادگیری ماشین بر اساس جریان و مبتنی بر ویژگی‌های بسته در مقالات پیشنهاد شده است، که ترافیک را به طور دقیق شناسایی می‌نماید. با این حال، کلاسه‌بندی ترافیک برای پروتکل‌های کپسوله شده (مانند استفاده از پروکسی سرور یا تونل‌های VPN) که عمدتاً برای مخفی کردن هویت کاربران به دلایل حفظ حریم خصوصی استفاده می‌شود، چالش برانگیز بوده و از این رو به طور گسترده مورد بررسی قرار نمی‌گیرند.

بعضی از این روش‌های مبتنی بر رویکردهای آماری و یادگیری ماشین با این فرض اولیه که ترافیک اساسی هر برنامه کاربردی برخی ویژگی‌های آماری دارد که تقریباً منحصر به هر برنامه است، عمدتاً به عنوان روش‌های آماری شناخته می‌شوند. در هر روش آماری از توابع و آمار آن استفاده می‌شود. بر اساس اطلاعات بدست آمده می‌توان گفت برای اولین بار در مرجع [۲] یک روش برای تشخیص ترافیک VPN به معنای وسیع و شناسایی ۷ نوع ترافیک مختلف پیشنهاد شده است.

تعداد زیادی مقاله با روش‌های یادگیری ماشین در زمینه کلاسه‌بندی ترافیک منتشر شده‌اند. در ادامه به بررسی دو مقاله بسیار مهمی ([۱] و [۲]) که بر روی این مجموعه داده بر اساس روش یادگیری ماشین منتشر شده است، می‌پردازیم. قبل از این مراجع، تنها یک مقاله بر اساس ایده‌های یادگیری عمیق گزارش شده است [۷]. که در آن از اتوانکدرهای انباشته<sup>۱</sup> (SAE) برای کلاسه‌بندی برخی از ترافیک‌های شبکه‌ای از پروتکل‌ها مانند HTTP و SMTP استفاده شده است. با این حال، در گزارش فنی آن، به مجموعه داده‌ی مورد استفاده اشاره‌ای نشده است. علاوه بر این، روش مورد استفاده در طرح، جزئیات پیاده‌سازی آن و گزارش مناسب نتایج نیز در آن وجود ندارد.

گیل و همکارانش [۲] از ویژگی‌های مرتبط با زمان مانند مدت زمان جریان، بایت‌های جریان بر ثانیه و زمان ورودی مسیر رفت و برگشت برای توصیف ترافیک شبکه با استفاده از الگوریتم‌های درخت تصمیم‌گیری k نزدیکترین همسایه (k-NN) و C<sub>۴,۵</sub> استفاده کرده‌اند. آنها تقریباً به ۹۲٪ دقت در فراخوانی<sup>۲</sup> دست یافته‌اند، که در آن مشخصه‌های اصلی ترافیک شامل مرورگر وب، ایمیل، چت، جریان، انتقال فایل و VoIP با استفاده از الگوریتم C<sub>۴,۵</sub> تعیین گشته است. همچنین آنها با استفاده از الگوریتم C<sub>۴,۵</sub> بر روی مجموعه داده‌هایی که از طریق تونل VPN می‌شوند، حدود ۸۸٪ دقت در فراخوانی را به دست آورده‌اند. اشکال اصلی این روش‌ها این است که استخراج ویژگی‌ها و مراحل انتخاب آنها اساساً با کمک یک متخصص صورت می‌پذیرد. از این رو، این کار باعث می‌شود که این رویکردها زمان‌گیر، گران و مستعد خطاهای انسانی باشد. علاوه بر این مدت زمانی که برای پیش‌بینی در الگوریتم k-NN استفاده می‌شود، نیز یک نگرانی عمده است.

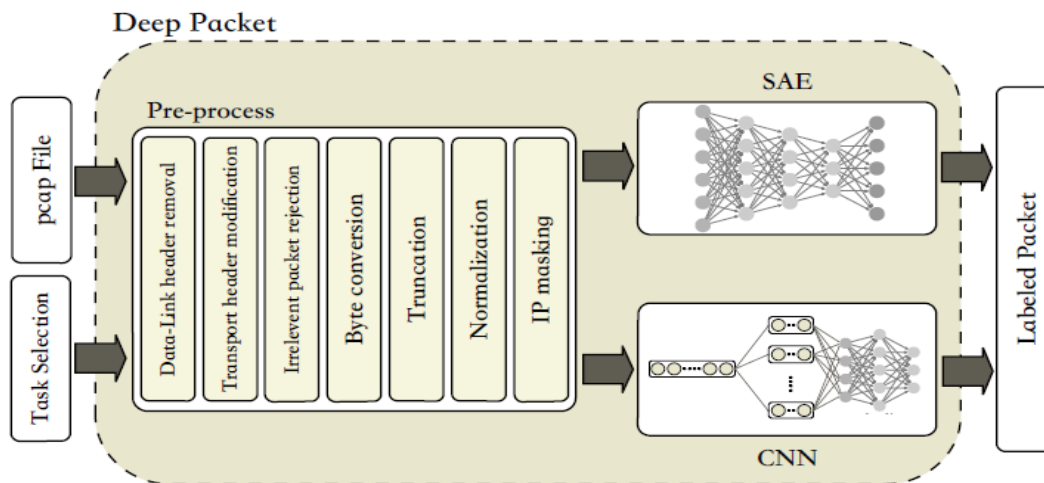
در مرجع [۱]، چارچوبی به نام بسته عمیق، که از دو روش یادگیری عمیق، یعنی شبکه‌های عصبی کانولوشنال<sup>۳</sup> و شبکه‌های عصبی اتوانکدر انباشته به وجود آمده که هم برای «شناسایی کاربرد» و هم «ویژگی‌های ترافیک» از

<sup>۱</sup> Stacked Auto-Encoder

<sup>۲</sup> recall

<sup>۳</sup> Convolutional NN

آن استفاده شده است. قبل از آموزش شبکه‌های عصبی، باید داده‌های ترافیک شبکه را تهیه کنیم تا بتوان به درستی آنها را به شبکه منتقل کرد. برای این منظور، مرحله پیش پردازش را در مجموعه داده انجام می‌دهیم. شکل ۱ ساختار کلی بسته عمیق را نشان می‌دهد. در مرحله تست، از یک شبکه عصبی از پیش آموزش دیده مطابق با نوع کلاسه‌بندی، شناسایی برنامه یا مشخصه ترافیک، برای پیش بینی کلاس ترافیکی که بسته به آن تعلق دارد، استفاده می‌شود. جزئیات مجموعه داده و معماری شبکه عصبی پیشنهادی [۱] در زیر شرح داده شده است.



شکل ۱: ساختار کلی بسته عمیق [۱]

## مدل پیشنهاد شده

در این بخش ابتدا مجموعه داده استفاده شده در این پروژه و مشخصاتش را معرفی کرده، سپس پیش پردازش‌های لازم برای آماده‌سازی آن را عنوان می‌نماییم. در نهایت نیز مدل پیشنهادی مرجع [۱] را بررسی کرده و مدل پیشنهادی خود را ارائه می‌دهیم.

## مجموعه داده و مشخصاتش [۲]

برای این پروژه، ما از مجموعه داده ترافیک ISCX VPN-nonVPN استفاده می‌کنیم که شامل ترافیک گرفته شده از برنامه‌های کاربردی مختلف در فایل‌هایی به فرمت pcap می‌باشد [۲]. در این مجموعه داده، بسته‌های گرفته شده به فایل‌های pcap مختلف بر اساس برنامه کاربردی تولیدکننده بسته‌ها (به عنوان مثال اسکایپ و Hangouts) و فعالیت خاصی که برنامه در مدت زمان ضبط بسته (مانند تماس صوتی، چت، انتقال فایل و یا تماس ویدیویی) انجام می‌دهد، تقسیم می‌شود.

برای ایجاد این مجموعه داده، ترافیک واقعی تولید شده توسط اعضای آزمایشگاه ضبط شده است. برای استفاده از خدماتی مانند اسکایپ و فیس بوک، حساب کاربری را برای دو کاربر ایجاد کرده‌اند. در جدول ۱ لیست کامل انواع مختلف ترافیک و برنامه‌های کاربردی موجود در مجموعه داده ارائه شده است. همانطور که از این جدول مشخص



است، در این مجموعه داده از ۱۷ نوع برنامه کاربردی استفاده شده و چون هر نوع ترافیک (VoIP، P2P و غیره) یک بار معمولی و بار دیگر با VPN ضبط شده، بنابراین مجموعاً ۱۲ دسته ترافیکی نیز داریم. ترافیک با استفاده از Wireshark و tcpdump ضبط شده و کل اطلاعات تولید شده ۲۸ گیگابایت است. برای ترافیک VPN، از یک ارائه‌دهنده خدمات VPN خارجی استفاده شده و با استفاده از OpenVPN به آن متصل گشته‌اند. برای ایجاد ترافیک SFTP و FTPS همچنین از یک ارائه‌دهنده خدمات خارجی و Filezilla به عنوان یک مشتری استفاده کرده‌اند. برای این پروژه این مجموعه داده را در طول چند روز دانلود نموده و همه مراحل را که در بخش بعد توضیح داده می‌شود بر آن اعمال کردیم.

جدول ۱: لیست انواع ترافیک و برنامه‌های کاربردی [۱]

Application	Size	Class Name	Size
AIM chat	5K	Chat	82K
Email	28K	Email	28K
Facebook	2502K	File Transfer	210K
FTPS	7872K	Streaming	1139K
Gmail	12K	Torrent	70K
Hangouts	3766K	VoIP	5120K
ICQ	7K	VPN: Chat	50K
Netflix	299K	VPN: File Transfer	251K
SCP	448K	VPN: Email	13K
SFTP	418K	VPN: Streaming	479K
Skype	2872K	VPN: Torrent	269K
Spotify	40K	VPN: VoIP	753K
Torrent	70K		
Tor	202K		
Voipbuster	842K		
Vimeo	146K		
YouTube	251K		

### پیش پردازش داده

از آنجایی که مجموعه داده از یک شبیه‌سازی واقعی بدست آمده است، لذا حاوی بسته‌های نامرتب نیز هست که باید حذف شوند. این بخش از داده برای انجام کارهایی مانند ایجاد یک اتصال یا اتمام آن مورد نیاز هستند، اما آنها هیچ اطلاعاتی در مورد برنامه کاربردی تولیدی آنها انتقال نمی‌دهند، پس می‌توان آنها را با خیال راحت از بین برد. علاوه بر این، بخش خدمات نام دامنه<sup>۱</sup>ی در این مجموعه داده وجود دارد، اطلاعاتی مربوط به شناسایی

<sup>۱</sup> Domain Name Service (DNS)

برنامه‌های کاربردی یا توضیحات تراکنشی را دربرنمی‌گیرند، از این رو می‌توان آن بخش را از مجموعه داده حذف کرد.

طول بسته در این مجموعه داده متغیر است، در حالی که استفاده از شبکه‌های عصبی در کل نیاز به ورودی با اندازه یکسان دارند. برای این منظور، قطع کردن بردار ورودی در یک طول مشخص یا اضافه کردن صفر<sup>۱</sup> به انتهای آنها اجتناب‌ناپذیر است. از مجموعه داده مشخص است که تقریباً ۹۶٪ بسته‌ها دارای بار<sup>۲</sup> با طول کمتر از ۱۴۸۰ بایت هستند. از این رو، ما هدر IP و اولین ۱۴۸۰ بایت بعد از آن از هر بسته IP را نگه می‌داریم بنابراین یک بردار ۱۵۰۰ بایتی به عنوان ورودی به شبکه‌های عصبی پیشنهادی اعمال می‌کنیم. برای بسته‌های دارای بارهای IP کمتر از ۱۴۸۰ بایت، در انتهای آنها صفر اضافه می‌شود. برای بدست آوردن عملکرد بهتر، تمام بایت‌های بسته را به ۲۵۵ یعنی حداکثر مقدار یک بایت تقسیم می‌نماییم، به طوری که تمام مقادیر ورودی در محدوده [۰ و ۱] قرار گیرد.

علاوه بر این، از آنجا که امکان دارد که شبکه عصبی در تلاش برای کلاسه‌بندی بسته‌ها از آدرس‌های IP موجود استفاده کند، و چون مجموعه داده با استفاده از تعداد محدودی میزبان<sup>۳</sup> و سرور ضبط شده است، پس برای جلوگیری از **Overfitting** آدرس‌های IP موجود در هدر IP را حذف می‌کنیم. با اعمال این مراحل پیش پردازش بر روی مجموعه داده می‌توان اطمینان داشت که شبکه عصبی هیچ استفاده‌ای از ویژگی‌های نامناسب برای انجام کلاسه‌بندی نخواهد کرد.

### معماری پیشنهادی

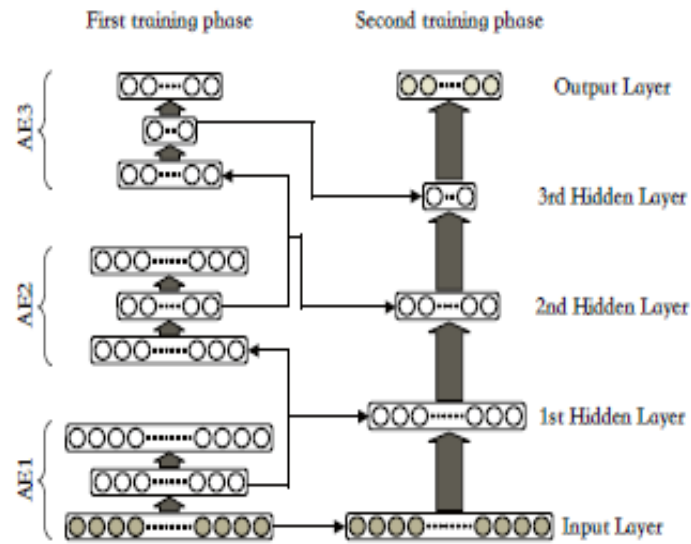
معماری SAE پیشنهادی [۱] شامل پنج لایه کاملاً متصل است که به ترتیب از بیش از ۴۰۰، ۳۰۰، ۲۰۰، ۱۰۰ و ۵۰ نورون تشکیل شده است. برای جلوگیری از مشکل **overfitting**، پس از هر لایه، تکنیک **dropout** با نرخ ۰.۰۵ مورد استفاده قرار می‌گیرد. در این تکنیک، در فاز آموزش، برخی از نورون‌ها به طور تصادفی صفر در نظر گرفته می‌شوند. بنابراین، در هر مرحله تکرار، یک مجموعه تصادفی از نورون‌های فعال وجود دارد. برای شناسایی برنامه‌های کاربردی و وظایف مشخصه ترافیک، در لایه نهایی SAE پیشنهادی، یک کلاسه‌بندی کننده از نوع **softmax** به ترتیب با ۱۷ و ۱۲ نورون افزوده می‌شود.

---

<sup>۱</sup> zero padding

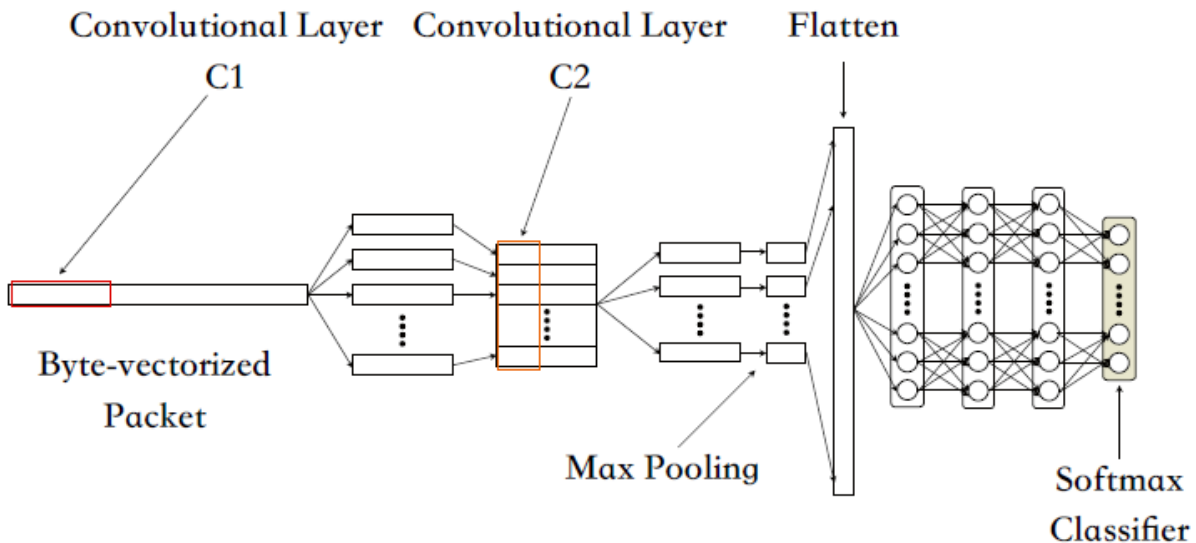
<sup>۲</sup> Payload

<sup>۳</sup> host



شکل ۲: آتوانکدر انباشته

تصویر کوچکتر از طرح پیشنهادی دوم [۱]، بر اساس CNN تک بعدی، در شکل ۳ نشان داده شده است. این مدل شامل دو لایه کانولوشن متوالی و به دنبال آن یک لایه پولینگ است. سپس تانسور دو بعدی به یک بردار یک بعدی تبدیل می‌شود و به شبکه سه لایه‌ای از نورون‌های کاملاً متصل منتقل می‌شود و همچنین از تکنیک dropout برای جلوگیری از overfitting استفاده می‌شود. در نهایت، یک کلاسه‌بندی کننده از نوع softmax به منظور کلاسه‌بندی، مشابه معماری SAE اعمال می‌شود. بهترین مقادیر برای لایه‌های کانولوشن در جدول ۲ نشان داده شده است.



شکل ۳: معماری CNN تک بعدی پیشنهادی [۱]

جدول ۲: بهترین مقادیر برای لایه‌های کانولوشن [۱]

Task	C1 Filter			C2 Filter		
	Size	Number	Stride	Size	Number	Stride
App. Idn.	4	200	3	5	200	1
Traffic Char.	5	200	3	4	200	3

در معماری پیشنهادی این پروژه، به جای CNN تک بعدی، شبکه کپسوله<sup>۱</sup> پیشنهاد می‌گردد. در این معماری به جای استفاده از دو لایه کانولوشن، یک لایه کپسول جاگذاری می‌شود که با توجه به ساختار شبکه‌های کپسوله دیگر نیازی به لایه پولینگ وجود ندارد. در این حالت نیز مانند مدل قبل از Dropout و در لایه آخر از Softmax استفاده می‌گردد. از مزایای این روش نسبت به مدل مرجع [۱] تعداد پارامترهای کمتر، سرعت همگرایی و دقت بیشتر است.

## نتایج

در مرجع [۱] نشان داده شده که شبکه عصبی کانولوشنال دارای نتایج بهتری هم در کلاسه‌بندی برنامه‌های کاربردی و هم در نوع ترافیک شبکه، نسبت به آتوانکدر انباشته می‌باشد. همچنین طبق جدول ۳ نسبت به مراجع قبلی که مدل خود را بر روی همین مجموعه داده پیاده‌سازی نموده‌اند نیز دارای نتایج بهتری است.

جدول ۳: مقایسه نتایج حاصل از روش بسته عمیق و مراجع دیگر

Paper	Task	Metric	Results	Algorithm
Deep Packet [47]	Application Identification	Accuracy	0.98	CNN
			0.94	k-NN
Deep Packet [16]	Traffic Characterization	Precision	0.93	CNN
			0.90	C4.5

حال برای بهبود هر چه بیشتر این نتایج فقط کافی است از شبکه‌ای مانند شبکه کپسوله که در کل ویژگی‌های مثبتی نسبت به شبکه‌های کانولوشنی دارد به جای آن استفاده نماییم. با توجه به توضیحات ارائه شده مبنی بر دشواری کار با فایل‌های با فرمت pcap، پیش پردازش آنها و پیاده‌سازی این حجم (۲۸ گیگا بایت) از اطلاعات بر روی سرورهای معمولی که آن را تقریباً امکان‌ناپذیر می‌کند، موفق به یاد دادن این شبکه جدید پیشنهادی نشدیم. بنابراین به نظر می‌رسد یکی از کارهای مهم می‌تواند این باشد که این مجموعه داده را به کل از فرمت pcap خارج

<sup>۱</sup> Capsule Networks

کرده و تبدیل به مجموعه داده دارای برچسب نرمالی با حجم کمتر نمود. مرحله بعدی آن نیز اعمال این مجموعه داده جدید به مدل پیشنهادی ما و مقایسه نتایج حاصل با نتایج مرجع [۱] می‌باشد.

## تحلیل

در این پروژه، همانند روش بسته عمیق مرجع [۱] روشی ارائه شده است که به طور خودکار ویژگی‌های ترافیک شبکه را از طریق الگوریتم‌های یادگیری عمیق برای کلاسه‌بندی ترافیک استخراج می‌کند. به همین علت، بسته عمیق اولین سیستم کلاسه‌بندی ترافیک با استفاده از الگوریتم‌های یادگیری عمیق، یعنی SAE و CNN تک بعدی است که می‌تواند هر دو وظیفه شناسایی برنامه‌های کاربردی و ویژگی‌های ترافیک را همزمان انجام دهد. بنابراین می‌توان روش بسته عمیق را گام اول در راستای یک روند کلی استفاده از الگوریتم‌های یادگیری عمیق برای کلاسه‌بندی ترافیک دانست. علاوه بر این، می‌توان با استفاده از این روش در هزینه استفاده از کارشناسان برای شناسایی و استخراج ویژگی‌های دستی از ترافیک صرفه‌جویی نمود و در نهایت منجر به کلاسه‌بندی دقیق‌تر ترافیک می‌شود. روش پیشنهادی در این پروژه را می‌توان نمونه تعمیم‌یافته بسته عمیق دانست که به جای CNN تک بعدی از شبکه کپسوله در آن استفاده شده است. که همه مزایای مربوط به بسته عمیق را به همراه مزایای شبکه‌های کپسوله دارا می‌باشد.

- [١] M. Lotfollahi, M. Jafari Siavoshani, R. Shirali Hossein Zade, M. Saberian, "Deep Packet: A Novel Approach For Encrypted Traffic Classification Using Deep Learning," *arXiv: ١٧٠٩.٠٢٦٥٦٧٣ [cs.LG]*, ٤ Jul ٢٠١٨.
- [٢] G. D. Gil, A. H. Lashkari, M. Mamun, A. A. Ghorbani, "Characterization of encrypted and vpn traffic using time-related features," in *٢nd International Conference on Information Systems Security and Privacy (ICISSP ٢٠١٦)*, ٢٠١٦.
- [٣] B. Yamansavascular, M. A. Guvensan, A. G. Yavuz, M. E. Karsligil, "Application identification via network traffic classification," in *Computing, Networking and Communications (ICNC), ٢٠١٧ International Conference on IEEE*, ٢٠١٧.
- [٤] V. Paxson, S. Floyd, "Wide area traffic: The failure of poisson modeling," *IEEE/ACM Transactions on Networking*, vol. ٣, no. ٣, pp. ٢٢٦-٢٤٤, ١٩٩٥.
- [٥] D. Wang, L. Zhang, Z. Yuan, Y. Xue, Y. Dong, "Characterizing application behaviors for classifying p٢p traffic," in *International Conference on Computing, Networking and Communications, ICNC, IEEE*, ٢٠١٤.
- [٦] J. Sherry, C. Lan, R. A. Popa, S. Ratnasamy, "Blindbox: Deep packet inspection over encrypted traffic," in *Proceedings of the ٢٠١٥ ACM Conference on Special Interest Group on Data Communication, SIGCOMM, ACM, New York, NY, USA*, ٢٠١٥.
- [٧] Z. Wang, "The applications of deep learning on traffic identification," *BlackHat USA*, ٢٠١٥.