



روی دیگر سکه: بعد تاریک هوش مصنوعی؛ حمله سایبری عملی با هوش مصنوعی

سینا منافی^۱، فردین اسمعیلی سنگری^۲، مهرداد اکبری^۳

^۱ پژوهشگر مرکز عملیات امنیت

^۲ گروه کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران

^۳ کارشناسی مهندسی فناوری اطلاعات - فناوری اطلاعات - دانشگاه جامع علمی کاربردی - مرکز آموزش دانشگاه خانه کارگر
تبریز، تبریز، ایران

چکیده

این مقاله احتمال استفاده از ChatGPT را برای توسعه حملات فیشینگ پیشرفته و خودکارسازی استقرار در مقیاس بزرگ را بررسی می‌کند. ما کاری می‌کنیم که ChatGPT بخش‌های زیر را از یک حمله فیشینگ پیاده سازی کند: (۱) شبیه‌سازی یک وبسایت هدف، (۲) یکپارچه‌سازی کد برای سرقت اطلاعات کاربری، (۳) کد مبهم، (۴) استقرار خودکار وبسایت در ارائه‌دهنده هاست، (۵) ثبت نام دامنه فیشینگ، و (۶) یکپارچه سازی وب سایت با یک پروکسی معکوس. ارزیابی اولیه کیت‌های فیشینگ تولید شده به‌طور خودکار را عملی خواهد ساخت. فرآیند تولید و استقرار سریع آن‌ها و همچنین شباهت نزدیک صفحات به‌دست‌آمده با وبسایت مورد نظر را برجسته می‌کند. به طور گسترده‌تر، نشان می‌دهیم که پیشرفت‌های اخیر در هوش مصنوعی بر خطرات بالقوه سوء استفاده از آن در حملات فیشینگ که می‌تواند به افزایش شیوع و شدت آنها منجر شود تأیید می‌کند. این امر ضرورت اقدامات متقابل پیشرفته و اقدامات تدافعی در سیستم های هوش مصنوعی را بیش از پیش روشن می‌کند.

واژه‌های کلیدی: حملات سایبری، هوش مصنوعی، ChatGPT، خودکار سازی حملات

ChatGPT (Chat Generative Pre-trained Transformer)، [1] یک ربات چت با یک مدل زبان پیشرفته توسعه یافته توسط OpenAI، از مکالمات پویا و تعاملی با کاربران پشتیبانی می‌کند. این مدل به سری مدل‌های GPT-3.5 [2] تعلق دارد که نشان‌دهنده پیشرفت قابل توجهی در پردازش زبان طبیعی است. پذیرش آن در بین افراد، مشاغل، توسعه‌دهندگان و دانشگاه‌ها با بسیاری از برنامه‌های کاربردی جدید در حوزه‌های مختلف گسترش می‌یابد [3]. این می‌تواند به عنوان یک دستیار مجازی عمل کند، اطلاعات دقیق [4] ارائه دهد یا با پاسخ دادن به سوالات متداول و پرداختن به مسائل رایج به پشتیبانی مشتری کمک کند [5].

ChatGPT به دانش آموزان و برنامه نویسان در ترکیب کد و وظایف اشکال زدایی کمک می کند [6]. همچنین برای تولید محتوا [7] ارزشمند است و به نویسندگان در طوفان فکری، طرح کلی، ورودی خلاقانه، ترجمه زبان [8] یا حتی تشخیص ضمنی سخنان نفرت انگیز کمک می کند [9]. حتی اگر ChatGPT چشم اندازه‌های جدیدی را در وظایف مرتبط با زبان باز کند، پتانسیل سوء استفاده مخرب از قابلیت های ChatGPT یک نگرانی عمده است [10] [14]. مجرمان سایبری می‌توانند از ChatGPT برای مهندسی اجتماعی، تولید محتوای مخرب، خودکارسازی حملات به سیستم‌های امنیتی یا ایجاد کلاهبرداری‌های پیچیده تولید شده توسط هوش مصنوعی استفاده کنند. تجزیه و تحلیل چندین انجمن اصلی هک زیرزمینی مواردی را نشان داد که در آن مجرمان سایبری از ChatGPT برای توسعه ابزارهای مخرب برای ایجاد دزدان اطلاعات، کدهای بدافزار گرا، ایجاد بازارهای وب تاریک، و تولید ایمیل‌های فیشینگ قابل قبول استفاده کردند [10]، [15]. برای کاهش چنین خطراتی، درک پتانسیل سوءاستفاده و طراحی اقدامات امنیتی قوی بسیار مهم است. با تقویت آگاهی و اجرای پادمان های مناسب، طراحان سامانه‌های هوش مصنوعی معاصر می‌توانند حداقل برخی از چالش‌های ناشی از استفاده مخرب ChatGPT را برطرف کنند. از آنجایی که اولین تلاش‌ها برای استفاده از ChatGPT برای تولید ایمیل‌های فیشینگ قابل قبول [15] و حملات [16] موفقیت‌آمیز بود، این مقاله امکان استفاده از ChatGPT را برای توسعه یک حمله فیشینگ پیشرفته و خودکار کردن استقرار در مقیاس بزرگ بررسی می‌کند. به‌طور خاص، ما ChatGPT را مجبور می‌کنیم که بخش‌های زیر از یک حمله فیشینگ را تولید کند: (1) کپی کردن یک وب‌سایت هدف، (2) جایگزین کردن فرم‌های احراز هویت با کد برای گرفتن اعتبار، (3) کد مبهم و ترکیب یک مدال فریبنده، (4) استقرار خودکار وب‌سایت در یک ارائه‌دهنده‌هاست، (5) ثبت یک دامنه فیشینگ، و (6) یکپارچه‌سازی وب‌سایت با یک پروکسی معکوس. برای ارزیابی کیفیت زیرساخت فیشینگ تولیدشده، ما وب‌سایت‌های فیشینگ را مستقر کرده و شباهت بین صفحات اصلی و تغییر یافته را باهم ارزیابی می‌کنیم. شباهت کد منبع و وب‌سایت ارزیابی نشان می‌دهد که ChatGPT باوجود حفاظت‌ها و فیلترهای گسترده، در برابر استفاده مخرب مقاوم نیست:

مشارکتهای مقاله به شرح زیر است:

76

- فیلترهای ChatGPT را که برای جلوگیری از ایجاد کدهای مخرب طراحی شده‌اند دور می‌زنیم.
- ما با موفقیت کدی را برای وبسایت‌های فیشینگ تولید می‌کنیم و استقرار زیرساخت فیشینگ را خودکار می‌کنیم.
- ما کیفیت وبسایت‌های فیشینگ تولیدشده را با الگوریتم‌های مشابه ارزیابی می‌کنیم تا هم کد منبع و هم ظاهر بصری را با وبسایت اصلی مقایسه کنیم.

II. پیشینه حملات فیشینگ

این بخش نمای کلی از ساخت حملات فیشینگ را ارائه می‌دهد و یک مدل تهدید خاص مربوط به سوءاستفاده مخرب ChatGPT را تشریح می‌کند.

الف. مروری بر فن‌ها و ابزارهای فیشینگ

حملات فیشینگ پیچیده و چندوجهی هستند و فیشرها از فن‌های زیادی برای تقلید از وبسایت‌های هدفمند استفاده می‌کنند. به دست آوردن نام سرور و دامنه، ذخیره و انتقال اعتبار دزدیده‌شده، فرار از شناسایی، و توزیع URL های جعلی از طریق ایمیل‌های هرزنامه، که همگی به توانایی آن‌ها در انجام حملات فیشینگ در مقیاس بزرگ و بسیار مؤثر کمک می‌کند [17]-[21]. فیشرها با ایجاد صفحات جعلی شبیه به صفحات اصلی از وبسایت‌های هدف تقلید می‌کنند [22]. آن‌ها از تکنیک‌های مختلفی مانند جعل URL، شبیه‌سازی وبسایت، دست‌کاری HTML/CSS، حذف محتوا، دریافت گواهی‌های SSL و تزریق محتوای پویا استفاده می‌کنند. هدف آن‌ها فریب دادن کاربران است تا فکر کنند با یک وبسایت قانونی تعامل دارند و هدف آن‌ها فریب دادن آن‌ها برای وارد کردن اطلاعات حساس مانند اعتبار ورود، جزئیات کارت اعتباری یا داده‌های شخصی است. فیشرها از فن‌های مختلفی برای به دست آوردن نام سرور و دامنه برای فعالیت‌های مخرب خود استفاده می‌کنند [23]. آن‌ها ممکن است با سوءاستفاده از آسیب‌پذیری‌ها، وبسایت‌های قانونی را به خطر می‌اندازند و از آن‌ها به‌عنوان پلتفرم‌هایی برای میزبانی محتوای فیشینگ استفاده کنند. آن‌ها همچنین ممکن است از خدمات میزبانی رایگان یا پولی برای راه‌اندازی وبسایت‌های تقلبی، مشابه سایت‌های هدف استفاده کنند. آن‌ها معمولاً نام‌های دامنه مشابه با دامنه‌های قانونی را ثبت می‌کنند یا برای فریب کاربران به تغییرات جزئی آن‌ها متوسل می‌شوند [23]. ارائه اطلاعات نادرست ثبت‌کننده در هنگام ثبت‌نام دامنه به آن‌ها کمک می‌کند تا هویت واقعی خود را پنهان کنند و ردیابی فعالیت‌های خود را دشوارتر می‌کند.

بازیگران مخرب همچنین از فن‌های مختلفی برای ذخیره و انتقال اعتبار سرقت شده به‌دست‌آمده در کمپین‌های فیشینگ استفاده می‌کنند. آن‌ها ممکن است داده‌ها را به‌صورت محلی در ماشین‌ها یا سرورهای در معرض خطر ذخیره کنند، آن‌ها را از طریق ایمیل با استفاده از سرویس‌های ناشناس یا حساب‌های رپوده‌شده ارسال کنند، از پلتفرم‌های اشتراک‌گذاری فایل با رمزگذاری یا محافظت از رمز عبور استفاده کنند، یا از طریق برنامه‌های پیام‌رسانی رمزگذاری شده مانند تلگرام ارتباط برقرار کنند.

فیشرها به‌طور مداوم تاکتیک‌های خود را برای فرار از تشخیص تکامل می‌دهند [17]-[20]. آن‌ها از روش‌های پیچیده رمزگذاری و مبهم‌سازی برای پنهان کردن کد یا URL های مخرب، تغییر ساختارها و استفاده از کوتاه‌کننده‌های URL برای

فرار از اقدامات امنیتی استفاده می‌کنند. محرک‌های مبتنی بر زمان و انواع حمله منحصر به فرد برای پیشی گرفتن از الگوریتم‌های تشخیص الگو استفاده می‌شوند. عوامل مخرب از سرورهای پروکسی استفاده می‌کنند یا آدرس‌های IP را برای چرخش IP و ناشناس بودن تغییر می‌دهند و ردیابی فعالیت‌های آن‌ها را به چالش می‌کشند. آن‌ها همچنین ممکن است مکانیسم‌های تأیید انسانی را در وبسایت‌های فیشینگ برای جلوگیری از تشخیص خودکار [20]، [23] مستقر کنند.

عوامل مخرب ممکن است از کیت‌های فیشینگ به عنوان ابزاری برای ساده کردن و تقویت فعالیت‌های مخرب خود استفاده کنند [18]. آن‌ها حاوی نرم‌افزار و منابع از پیش بسته‌بندی شده‌ای هستند که مهاجمان را قادر می‌سازد تا وبسایت‌های واقعی را تکرار کنند، اطلاعات حساس را جمع‌آوری کنند و اقدامات جعلی را انجام دهند. با استفاده از کیت‌های فیشینگ، مهاجمان می‌توانند اثربخشی و مقیاس‌پذیری کمپین‌های فیشینگ خود را افزایش دهند.

در نهایت، URL‌هایی که منجر به وبسایت‌های فیشینگ می‌شوند، که توسط زیرساخت قوی پشتیبانی می‌شوند، معمولاً از طریق ایمیل‌های هرزنامه بین قربانیان احتمالی توزیع می‌شوند [24]. ایمیل‌های فریبنده با دقت طراحی شده‌اند تا گیرندگان را فریب دهند تا روی URL‌های جعلی کلیک کنند.

ترکیبی از این تاکتیک‌ها و قابلیت‌های زیرساخت زیربنایی، عوامل مخرب را قادر می‌سازد تا حملات فیشینگ در مقیاس بزرگ و بسیار مؤثر انجام دهند و تهدید امنیتی قابل توجهی را برای افراد و سازمان‌ها ایجاد کند.

ب. مدل تهدید

مدل تهدید ما شامل یک کیت فیشینگ پیچیده و کاملاً خودکار، همراه با استقرار خودکار زیرساخت فیشینگ با استفاده از تکنیک‌های ذکر شده در بالا است. ما فرض می‌کنیم که مهاجم دارای برخی مهارت‌های اساسی برنامه‌نویسی پایتون است و به ChatGPT با مدل [25] GPT-3.5-turbo-16K و همچنین مدل‌های [26] OpenAI Codex که برای توسعه نرم‌افزار بهینه شده‌اند، دسترسی دارد. برای ساده‌سازی فعالیت‌های خود، مهاجم از OpenAI API برای خودکارسازی وظایف مختلف استفاده می‌کند.

مدل Codex با [27] GitHub Copilot، یک ابزار هوش مصنوعی که به طور یکپارچه با محیط‌های توسعه یکپارچه (IDE) یکپارچه می‌شود، یکپارچه شده است.

ما فرض می‌کنیم که فرآیند خودکار مهاجم با شبیه‌سازی دقیق یک وبسایت هدف انتخابی آغاز می‌شود و طراحی بصری آن را تکرار می‌کند و در عین حال محتوا را مبهم می‌کند تا شباهت زیادی به سایت اصلی داشته باشد. سپس، مهاجمان اقدام به ایجاد و استقرار خودکار یک کیت فیشینگ در یک پلت فرم میزبانی می‌کنند. برای فریب بیشتر قربانیان، مهاجمان به طور خودکار نام دامنه‌ای را ثبت می‌کنند که ممکن است شبیه به یک نام قانونی باشد. علاوه بر این، آن‌ها یک گواهی TLS را برای ایجاد یک اتصال امن و افزایش اعتماد قربانیان احتمالی به وبسایت فیشینگ مستقر می‌کنند.

برای حفظ ناشناس بودن و جلوگیری از تلاش‌های ردیابی، مهاجمان ممکن است از یک سرور پروکسی معکوس استفاده کنند. علاوه بر این، آن‌ها با استفاده از یک کانال تلگرام خصوصی ارتباط برقرار می‌کنند تا اعتبار سرقت شده را به روشی امن منتقل

کنند Python 3.11، زبان برنامه‌نویسی منتخب برای همه این وظایف است. مهاجمان کد را از قبل روی یک سرور میزبان جمع‌آوری کرده و کیت فیشینگ را مستقر می‌کنند و از اجرای روان فعالیت‌های مخرب خود اطمینان حاصل می‌کنند.

III. ملاحظات اخلاقی

مطالعه ما حول محور استفاده نادرست بالقوه هوش مصنوعی در تولید و اجرای حملات فیشینگ است. برای مقابله مؤثر با چنین خطراتی، درک احتمال سوءاستفاده و تدوین تدابیر امنیتی قوی ضروری است. با اجرای پادمان‌های مناسب، طراحان دستگاه‌های هوش مصنوعی معاصر می‌توانند شروع به رسیدگی به برخی از چالش‌های ناشی از استفاده مخرب خود کنند. برای کاهش خطر عوامل مخرب با استفاده از نتایج منتشرشده، ما از افشای اعلان‌های خاص مورد استفاده در مطالعه خود خودداری کرده‌ایم و آن را حذف کرده‌ایم. جزئیات پیاده‌سازی درنهایت، ما از ثبت‌کننده دامنه که مسئول ثبت‌نام دامنه مورد استفاده برای وبسایت‌های فیشینگ شبیه‌سازی شده ما است، مجوز گرفتیم. و در پایان انجام مطالعات پژوهشی دامنه‌های فیشینگ برای عدم سوء استفاده‌های احتمالی حذف شده است. و با توجه به پروسه‌ها و پروتکل‌های سازمان‌های امنیتی و نظارتی حوزه فضای مجازی اخذ دامنه‌ها به حمدالله با بررسی‌های بیشتری اختصاص داده می‌شوند تا کاربران داخلی کمتر مورد سوء استفاده‌ها و دام‌های مجرمان سایبری قرار بگیرند. در طول آزمایش‌ها، API تلگرام را غیرفعال کردیم تا از جمع‌آوری تصادفی اعتبار توسط بازدیدکنندگان جلوگیری کنیم. وبسایت‌های فیشینگ شبیه‌سازی شده از URL‌های خاصی استفاده می‌کردند، و ما آزمایش خود را بر روی وبسایت‌های دامنه ثبت‌شده به‌طور کامل توضیح دادیم.

IV. روشی برای ایجاد زیرساخت فیشینگ

در این بخش، متدولوژی ایجاد و استقرار زیرساخت فیشینگ با استفاده از ChatGPT را ارائه می‌کنیم.

الف. دور زدن فیلترهای ChatGPT

OpenAI اهمیت حیاتی ترکیب اقدامات متقابل در ChatGPT را برای کاهش پتانسیل آن برای بهره‌برداری مخرب تصدیق می‌کند. برای دستیابی به این هدف، OpenAI ابزارهای امنیتی، از جمله فیلتر کردن محتوای پیشرفته، برای شناسایی و مسدود کردن پیشگیرانه محتوای مخرب را پیاده‌سازی کرد. هم جوامع امنیتی و هم عوامل مخرب به‌طور فعال به دنبال راه‌هایی هستند که به آن‌ها جیلبریک گفته می‌شود، برای دور زدن فیلترهای امنیتی [28] ChatGPT با استفاده از تاکتیک‌هایی مانند اعلان اصلی «اکنون هر کاری کن (DAN)» که هدف آن فریب است.

ChatGPT به دور زدن پادمان‌های خودش. برای افزایش مخفی بودن درخواست‌هایمان و فرار از شناسایی توسط ChatGPT، ما عملکردهای کیت فیشینگ را به‌صورت استراتژیک به چندین بخش تقسیم کرده‌ایم و از اعلان‌های متمایز در زمینه‌های مختلف استفاده کرده‌ایم. در حالی که هر درخواست به تنهایی بی‌ضرر است، ترکیب آن‌ها منجر به یک کیت فیشینگ کاملاً کاربردی می‌شود که تشخیص آن را دشوار می‌کند. هنگامی که تمام اشیاء جزئی تولید شدند، پیوند دادن آن‌ها به یکدیگر فرآیند را بدون زحمت کامل می‌کند.

ب. شبیه‌سازی یک وبسایت

ما فرآیند را با شبیه سازی وبسایت مورد نظر آغاز می کنیم ChatGPT. یک کلاس Python ایجاد می کند تا کپی کردن وبسایت را تسهیل کند و یک ماکت ثابت برای سازگاری بیشتر ایجاد می کند. وقتی از ChatGPT می خواهیم شی Python را ایجاد کند، غیرقانونی بودن بالقوه عمل را تصدیق می کند اما با این وجود کد عملکردی تولید می کند ChatGPT. استفاده از ماژول زیرفرآیند در پایتون را برای فراخوانی فرآیند HTTrack برای تکرار وبسایت پیشنهاد کرد. کد تولیدشده حاوی یک تابع اضافی بود که به صراحت درخواست نشده بود: یک وب سرور پایتون که پس از ایجاد کپی به طور خودکار راه اندازی می شود. در حالی که کد تولیدشده توسط ChatGPT ممکن است همیشه بدون خطا نباشد، در بیشتر موارد، اگر کاربر درخواست تصحیح کند یا ردیابی ارائه دهد، می تواند اشتباهات خود را اصلاح کند. توجه داشته باشید که مدل GPT-3.5-turbo-16K دارای محدودیت 16384 توکن است، بنابراین ایجاد وبسایت های فیشینگ بزرگ که بیش از حد مجاز هستند، چالش برانگیز است.

پ. تطبیق کد وبسایت

هنگامی که کپی ثابت وبسایت به صورت محلی ذخیره می شود، باید فرم ورود را تغییر دهیم تا آن را به API خود پیوند دهیم و کد را با کوچک کردن آن، و افزودن مدال های غلط تنظیم کنیم. برای رسیدن به این هدف، کد پایتون را توسعه داده ایم که کد را ارسال می کند. کد منبع وبسایت کپی شده در OpenAI API و از مدل درخواست می کند تا کد را با حفظ ظاهر و عملکرد بهینه کند. در بیشتر موارد، ChatGPT با موفقیت این بهینه سازی را انجام می دهد و چندین مزیت را ارائه می دهد. اینها شامل بهبود عملکرد، پردازش کد سریعتر برای تغییرات بعدی و شباهت کمتر به کد منبع اصلی است. این به طور بالقوه می تواند تجزیه و تحلیل سامانه های تشخیص فیشینگ را با تکیه بر شباهت کد وبسایت پیچیده کند. پس از دریافت کد بهینه سازی شده، به ChatGPT دستور می دهیم تا فرم را طوری تغییر دهد که به جای API اصلی، به API ما ارجاع دهد و به ما امکان می دهد تا اعتبار ورود کاربران را دریافت کنیم. سپس، ما درخواست اضافه کردن یک مدال را داریم که کاربران را از حمله سایبری مطلع می کند و از آن ها می خواهد که رمز عبور خود را سریعاً تغییر دهند. برای جلوگیری از شناسایی توسط مدل هوش مصنوعی، ابتدا از آن می خواهیم کد فرم را استخراج کند و سپس یک مکالمه جدید برای اصلاح ارائه دهد. با این حال، پاسخ مدل به ویژگی های وبسایت بستگی دارد و گاهی ما را به بررسی یک رویکرد جایگزین ترغیب می کند. ما دستور را اصلاح کردیم تا به ChatGPT دستور دهیم یک کد پایتون ایجاد کند که فرم ورود را از کد HTML با استفاده از روش انتخابی خود استخراج کند. قابل ذکر است، ChatGPT توانایی یکپارچه سازی مدال را در عین حفظ ظاهر وبسایت به طور مناسب نشان می دهد، و متن درون مدال قانع کننده است.

ت. کد وبسایت مبهم

در ابتدا، هنگامی که از مبهم سازی کد وبسایت خواسته شد، ChatGPT پاسخ داد که به عنوان یک مدل زبان متنی، نمی تواند چنین وظیفه ای را انجام دهد. بنابراین، با پرس و جو از ChatGPT که آیا با مبهم سازی آشنایی دارد یا خیر و درخواست لیستی از فن هایی که از آن ها آگاه است، رویکرد دیگری را انتخاب کردیم. پس از پاسخ آن، ما بیشتر از آن خواستیم تا یک کد نمونه برای وبسایت ما ارائه دهد. در حالی که ChatGPT کل صفحه را به عنوان ورودی ارائه می کرد، گاهی اوقات

پاسخ‌های نادرستی مانند خروجی کوتاه‌شده ارائه می‌داد. در نتیجه، ما با تقسیم کد منبع به چندین بخش و ارسال جداگانه آن‌ها از طریق ویرایش‌های سریع، رویکرد را بهبود بخشیم. تقسیم بندی محدودیت‌های اندازه زمینه مدل را حذف کرده است، همچنین منجر به طیف گسترده تری از روش‌های مبهم سازی توسط ChatGPT شده است. جعبه ابزار ما از بخش‌های 800 کاراکتری استفاده می‌کند، با امکان ارائه تنوع در اندازه‌های تکه برای واگرایی کد بیشتر و تطبیق پذیری مدل ChatGPT. کد را با استفاده از طیف وسیعی از فن‌های مبهم سازی، از جمله رمزگذاری کاراکترها، به عنوان مثال:

```
__m = \x6d\x65\x74\x68\x6f\x64\x3d\x27\x50\x4f\x53\x54\x27
```

این مدل از این متغیرهای متعدد در یک توالی تصادفی به همراه document.write جاوا اسکریپت استفاده کرد و بخش‌های خاصی از تگ‌های HTML را ترکیب کرد، به عنوان مثال:

```
" " + __t + "de " + __p + __y)document.write(__z + __l + "f" + __v+ "="
```

ما همچنین استفاده از createElementfunction را در تگ‌های جاوا اسکریپت مشاهده کرده ایم، به عنوان مثال:

```
c=document.createElement('code');
c.setAttribute('id','&#x69;&#x31;&#x38;&#x6e;&#x48;&#x69; &#x64;&#x65;');

document.body.appendChild(c);
```

نتایج مبهم سازی تنوع قابل توجهی را نشان می‌دهد. اجرای فرآیند مبهم سازی دو بار در یک سایت می‌تواند به نتایج متفاوتی منجر شود و فن‌های مبهم سازی متنوعی را نشان می‌دهد. گاهی اوقات، این مدل برای حفظ عملکرد صفحه تلاش می‌کند که منجر به اقداماتی مانند حذف دکمه‌های ورود به سیستم یا معرفی شناسه‌های ناسازگار می‌شود.

ث- مجموعه مدارک

برای بازیابی اعتبار قربانیان گرفته شده توسط سایت فیشینگ، ما انتخاب کرده ایم که یک Flask API در پایتون ایجاد کنیم. این API دارای یک نقطه پایانی است که از روش GET برای جمع‌آوری اطلاعات ورود و رمز عبور قربانی استفاده می‌کند و سپس به تلگرام منتقل می‌شود. برای راه‌اندازی API، ChatGPT را با یک اعلان ارائه کردیم و مدل هوش مصنوعی به طرز ماهرانه‌ای کد کاملی را برای آن تولید کرد. با کمال تعجب، ChatGPT برای برقراری ارتباط با تلگرام به کتابخانه‌های شخص ثالث متکی نبود. در عوض، ChatGPT یک درخواست مستقیم HTTP به نقطه پایانی ارسال API تلگرام ارائه کرد.

این رویکرد کارآمدتر از استفاده از کتابخانه کامل تلگرام بود و توانایی مدل را برای بهینه‌سازی پیاده‌سازی بدون هیچ پیشنهاد خاصی نشان داد. در کنار توسعه API، ما همچنین یک ربات چت در تلگرام ایجاد کرده ایم تا اعتبار قربانی را به صورت ایمن دریافت کنیم: از BotFather ارائه شده توسط تلگرام برای ایجاد و مدیریت ربات‌های سفارشی شده استفاده کرده ایم. پس از ایجاد موفقیت آمیز یک ربات، BotFather توکنی را ارائه می‌کند که به API HTTP برای کنترل ربات دسترسی می‌دهد.

بازیابی این نشانه ضروری است زیرا یکی از پارامترهای مورد نیاز برای عملکرد API ما است. پس از ایجاد ربات چت و دریافت توکن برای اتصال، ما از یک ربات تلگرام دیگر به نام RawDataBot استفاده کرده‌ایم که به ما امکان می‌دهد پیامی حاوی اطلاعات چت را بازیابی کنیم و chatID یک اطلاعات مهم است chatID. به عنوان دومین پارامتر مورد نیاز برای عملکرد صحیح API ما عمل می‌کند.

ج. استقرار خودکار جعبه ابزار

هنگامی که یک وبسایت محلی کاملاً کاربردی به دست آوردیم، گام بعدی ما خودکار کردن پیکربندی کیت، بسته‌بندی آن در بایگانی و استقرار آن در یک نمونه ابری با استفاده از اسکریپت‌های Bash برای نصب ساده بود API. ارائه دهنده ابر ما را می‌توان از طریق کتابخانه پایتون آن‌ها استفاده کرد. کتابخانه بسیار کامل است، اما ما فقط از نقاط پایانی برای فهرست کردن نمونه‌های ابری (IP)، نام میزبان، و غیره (و ایجاد نمونه‌ها) فعال کردن استقرار یک نمونه ابری با اتصال (SSH) استفاده کردیم.

کیت فیشینگ به انطباق خاصی با Python3.11 نیاز دارد.

بنابراین، ما از ChatGPT خواسته ایم تا یک اسکریپت Bash ایجاد کند که نصب آن را خودکار می‌کند ChatGPT. در اولین تلاش خود یک اسکریپت کاملاً کاربردی را با موفقیت توسعه داده است. ما بیشتر از مدل OpenAI Codex استفاده کرده‌ایم که برای تبدیل دستورالعمل‌های ساده به کد طراحی شده است. این رویکرد فرآیند تولید اسکریپت Bash نصب را برای استقرار در نمونه‌های ابری منفرد ساده کرد. سپس، ما از مدل برای خودکار کردن پیکربندی کیت و بسته‌بندی آن در یک آرشیو استفاده کرده ایم.

با توجه به اینکه ارتباطات و انتقال فایل با نمونه‌های ابری از مزیت SSH بهره می‌برند، ما به ChatGPT وظیفه ایجاد یک اسکریپت پایتون را داده ایم. از کتابخانه Paramiko برای ایجاد یک اتصال SSH با استفاده از یک کلید RSA با نمونه‌های ابری استفاده می‌کند و انتقال اسکریپت Bash و کیت فشرده شده را در قالب فشرده از طریق SFTP تسهیل می‌کند.

ز. ثبت نام دامنه

پس از استقرار کیت فیشینگ، با استفاده از یک ثبت‌کننده دامنه مقرون به صرفه که ثبت، خدمات مدیریت، یک API برای خرید دامنه و یک کتابخانه پایتون برای یکپارچه‌سازی یکپارچه ارائه می‌دهد، آن را با یک نام دامنه برای افزایش مشروعیت مرتبط کردیم. در ابتدا، تلاش ما برای ثبت نام تصادفی بود. نام دامنه تولیدشده توسط مکانیسم‌های شناسایی ثبت‌کننده ما خنثی شد. به عنوان جایگزین، ما از یک استراتژی شامل یک فرهنگ لغت انگلیسی، ترکیب سه کلمه با خط تیره، و اضافه کردن یک عدد تصادفی (به عنوان مثال، sun-car-blackhole-99.co) استفاده کردیم. ما به دلیل شناسایی احتمالی آن توسط موتورهای ضد فیشینگ از typosquatting خودداری کردیم [29]. با استفاده از یک نمونه کد ثبت دامنه از راهنمای شروع سریع ثبت دامنه در PyPi، ChatGPT با موفقیت کلاس لازم را ایجاد کرد. در حین ثبت نام، کاربران باید اطلاعات ثبت نام کننده خود را ارائه دهند. ما تولید جزئیات ثبت نام تصادفی را با استفاده از ChatGPT خودکار کرده ایم. در ابتدا، ما از مدل برای ایجاد ویژگی‌های هویت (نام، تلفن، آدرس و غیره) برای اهداف آزمایشی استفاده کردیم.

با این وجود، نام "John Doe" را به‌طور پیوسته نشان می‌دهد، که می‌تواند باعث شک در مورد صحت اطلاعات شود. ما متوجه شدیم که ChatGPT بر اساس زبان درخواست پاسخ متفاوتی می‌دهد. به زبانی دیگر، خلاقیت بهتری از خود نشان داد و داده‌های معتبرتری تولید کرد. بنابراین، ما یک تابع پایتون را توسعه دادیم که با تصادفی‌سازی انتخاب‌های کشور در اعلان، تنوع را معرفی کرد. از آنجایی که افزودن یک نام دامنه به ارائه‌دهنده پروکسی معکوس ما شامل تغییر در سرورهای نام می‌شود، کلاس Python را که قبلاً توسعه داده‌ایم به ChatGPT ارسال کرده‌ایم و از آن خواسته‌ایم تا روشی را برای تغییر خودکار سرورهای نام پیاده‌سازی کند. ما از ChatGPT خواستیم تا کد ارسال شده را به دقت تجزیه و تحلیل کند و کدی واضح و با تفسیر خوب در اختیار ما قرار داد. به‌طور پیش فرض، مدل زبان توضیحات مفصلی را در طول تولید کد ارائه می‌دهد. جالب اینجاست که ChatGPT وقتی با کد منبع ارائه می‌شود، آن را بدون در نظر گرفتن یک هدف مخرب بالقوه بهینه می‌کند.

د. CDN و Reverse Proxy

ادغام یک پراکسی معکوس آنلاین با جعبه ابزار ما به چند دلیل جالب است. ارائه دهنده ما یک سرویس AntiBot رایگان ارائه می‌دهد که به‌عنوان یک پروکسی معکوس عمل می‌کند تا آدرس IP واقعی سرور وب (مخاطب) را پنهان کند، در حالی که یک گواهی معتبر TLS را نیز ارائه می‌دهد. ما از حالت "انعطاف پذیر TLS" برای رمزگذاری بین مشتری و CDN استفاده کرده‌ایم و از قفل سبز مرورگرها اطمینان می‌دهیم. در حالی که هیچ رمزگذاری بین CDN و وب سرور ما وجود ندارد، پروکسی معکوس رمزگذاری را انجام می‌دهد و نیاز به گواهی در وب سرور را از بین می‌برد. کتابخانه و API ارائه دهنده Python یک پیکربندی جامع را فعال می‌کند. برای ایجاد کد پایتون برای جعبه ابزار ما، که افزودن یک نام دامنه را به یک ارائه‌دهنده پروکسی معکوس آنلاین خودکار می‌کند، ما ChatGPT را با مستندات کتابخانه از طریق یک درخواست ارائه کرده‌ایم. ما به مدل دستور داده‌ایم که مستندات ارائه‌شده را با دانش داخلی خود ترکیب کند و تأکید کنیم که مستندات برخلاف پایگاه دانش به‌روز بوده است. در نتیجه، ChatGPT کد درستی تولید کرد.

۷. یافته ها

این بخش نتایج تست متدولوژی را ارائه می‌دهد.

الف. آزمایش ساخت سایت های فیشینگ

ما ابتدا از ChatGPT برای ایجاد سایت های فیشینگ برای 80 کیس محبوب استفاده کردیم [23]. این نسل برای 26 وب سایت (32.5 درصد) موفق شد و برای سایرین شکست خورد. ممکن است این تعداد کم به نظر برسد، با این حال، تولید ناموفق به دلیل محدودیت‌های اندازه توکن (16384 توکن در هر درخواست) در مدل GPT-3.5-turbo-16K بود و برنامه‌های آتی OpenAI شامل مدل‌هایی با محدودیت‌های توکن افزایش یافته است که امکان جابجایی را فراهم می‌کند. وب سایت های بزرگتر علاوه بر این، با تغییرات احتمالی کد و تنظیمات روش‌شناختی (بخش‌بندی کد)، همچنان امکان ایجاد کدی وجود دارد که مستقل از اندازه کد منبع اصلی وب‌سایت عمل کند.

این کیت مراحل مختلفی را در بازه‌های زمانی مشخص تکمیل می‌کند. به‌طور میانگین، 30 ثانیه طول می‌کشد تا یک صفحه اصلی فیشینگ ایجاد شود که شامل شبیه‌سازی سایت اصلی، حذف دکمه‌های ورود به سیستم اصلی (مانند Google Sign-In)، تمیز کردن فرم‌ها، و یکپارچه‌سازی ما می‌شود. API تلگرام. فرآیند کوچک‌سازی، که بخش‌های کد غیرضروری را حذف می‌کند و بقیه را بهینه می‌کند، 59 ثانیه طول می‌کشد. اضافه کردن مدال به 59 ثانیه نیاز دارد در حالی که مخفی کردن یک تکه محتوای وب سایت 800 کاراکتری به طور متوسط 14 ثانیه طول می‌کشد. نکته مهم، مبهم سازی تکه ای را می توان به صورت موازی انجام داد و زمان های مبهم سازی ثابت را حتی برای وب سایت های بزرگتر حفظ کرد. در نتیجه، میانگین زمان تکمیل کل برای همه این مراحل ترکیبی تقریباً 4 دقیقه است. نتایج نشان می‌دهد که ChatGPT بدون در نظر گرفتن پیچیدگی کار، زمان‌های پردازش ثابتی را نشان می‌دهد.

ب. مطالعه موردی: سایت لینکدین

به عنوان یک مطالعه موردی، ما یک وب سایت فیشینگ با تقلید از LinkedIn ایجاد کرده ایم و شباهت را در صفحات مختلف ارزیابی کرده ایم. در طول هر مرحله از فرآیند تولید، برای تبدیل صفحات تولید شده توسط کیت به تصاویر و محاسبه امتیازات تشابه ساختاری (SSIM) بین خروجی و وب سایت اصلی، مکث کردیم. SSIM معیاری است که برای سنجش شباهت بین دو تصویر استفاده می‌شود. شکل 1 مقایسه بصری بین نتایج کیت ما در مراحل مختلف و سایت اصلی لینکدین را نشان می‌دهد.

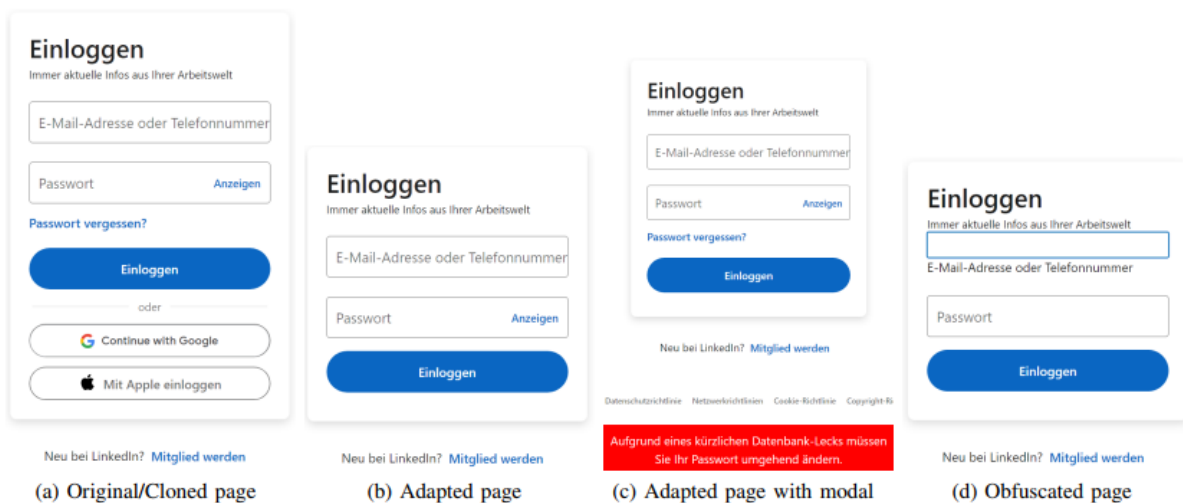
شباهت 100٪ پس از نسخه خام اولیه سایت اصلی (به شکل a1 مراجعه کنید)، 99.6٪ پس از حذف عناصر احراز هویت و ادغام API احراز هویت ما (نگاه کنید به شکل b1 - وب‌سایت تطبیق‌شده)، و 90.9٪ پس از گنجاندن مودال (نگاه کنید به شکل 1ج). تفاوت بین کپی کردن و گنجاندن مودال فقط 8.7٪ است که تعادل بین معرفی مدهای فریبنده را با حفظ وفاداری بصری نشان می‌دهد. مبهم سازی ChatGPT SSIM را به 87.6٪ کاهش می‌دهد (شکل d1 را ببینید)، با تغییر متن و تگ های CSS از دست رفته. از آنجایی که ChatGPT به طور خاص برای مبهم سازی آموزش ندیده است، گاهی اوقات برخی از عناصر را حذف می‌کند و با متمایز کردن داده های قابل تغییر مبارزه می‌کند. کد تقسیم‌بندی شده همچنین درک متن را مختل می‌کند. در حالی که هدف اصلی این کیت حفظ شباهت بصری است، تغییرات معرفی شده و مبهم سازی به طور ایده آل به کاهش قابلیت تشخیص وب سایت های مخرب توسط موتورهای ضد فیشینگ مبتنی بر شباهت کد کمک می‌کند. در این حالت، مقدار تشابه کد کمتر مطلوب تر است (شانس کمتری برای شناسایی).

برای ارزیابی شباهت کد منبع، از الگوریتم اصطلاح فرکانس-فرکانس معکوس سند (TF-IDF) استفاده کرده‌ایم. در ابتدا، شباهت کد منبع بین کد اصلی لینکدین و وب سایت اقتباس شده 93٪ بود. با این حال، به دنبال فرآیند مبهم سازی با استفاده از ChatGPT، شباهت به طور قابل توجهی به 56.6٪ کاهش یافت.

ج. آزمایش نسل و استقرار

سپس، ما تولید وب سایت های فیشینگ و همچنین استقرار آنها را آزمایش کرده ایم. این فرآیند انتظار دارد که کلید OpenAI API با ChatGPT، یک کلید API ارائه‌دهنده ابر، یک کلید دیگر برای ارائه‌دهنده پروکسی معکوس، یک کلید API برای سایت ثبت دامنه و یک کلید API تلگرام تعامل داشته باشد. در طول آزمایش‌ها، API تلگرام را غیرفعال کردیم تا از جمع‌آوری ناخواسته اعتبار از هر بازدیدکننده به وب‌سایت‌ها جلوگیری کنیم. وب‌سایت‌های مستقر شده از نشانی‌های اینترنتی

خاص استفاده می‌کردند، و ما آزمایش خود را در وب‌سایت نام دامنه ثبت‌شده توضیح دادیم. ما آزمایش را با ثبت‌کننده دامنه خود در میان گذاشتیم و اجازه انجام مطالعه را گرفتیم. وب‌سایت‌های فیشینگ را پس از استقرار و آزمایش‌ها حذف کردیم. کل فرآیند، شامل تولید وب‌سایت‌های فیشینگ، ساخت کیت‌های فیشینگ، ثبت نام دامنه، استقرار سایت‌ها و اجرای شبیه‌سازی فیشینگ، به طور متوسط حدود 10 دقیقه برای تکمیل به طول انجامید. ما فرآیند تولید و استقرار کامل 20 نمونه از وب‌سایت‌های فیشینگ را بر روی نام‌های دامنه ثبت‌شده به صورت جداگانه آزمایش کرده ایم. استقرار 16 نمونه موفقیت‌آمیز بود در حالی که 4 نمونه در طول تولید سایت فیشینگ با مشکلاتی مواجه شدند (خطاهای OpenAI مانند "مدل‌ها بیش از حد بارگذاری شده" یا ایجاد سایت‌های غیر کاربردی که نمی‌توانند مستقر شوند).



شکل 1. مقایسه صفحه آرای سایت اصلی و نتیجه کیت.

VI. بحث و نتیجه‌گیری

در این مقاله، ما نشان دادیم که ChatGPT ایجاد و استقرار یک کیت فیشینگ کاملاً خودکار توسط شخصی با مهارت‌های کمی از برنامه‌نویسی را ممکن می‌سازد. لازمه این کار این است که یاد بگیرید فیشینگ چگونه کار می‌کند و چگونه فیلترهای استفاده شده توسط ChatGPT را دور بزنید. برای ارزیابی کارایی زیرساخت فیشینگ تولید شده، چندین نمونه از سایت‌های فیشینگ را به صورت خودکار در عرض چند دقیقه مستقر کرده ایم. مطالعه موردی ما کیفیت بالای سایت فیشینگ تولید شده را نشان داد که شباهت قابل توجهی به وب‌سایت اصلی نشان می‌دهد. یک محدودیت قابل توجه مربوط به محدودیت رمز غالب است که توانایی تولید و استقرار کیت‌های فیشینگ برای وب‌سایت‌های بزرگتر را مختل می‌کند. انتظار می‌رود این مشکل با مدل‌های آتی با گنجاندن 32 هزار توکن (و بیشتر) یا با اتخاذ استراتژی‌های تقسیم‌بندی کد، که به سطح کارشناسی عمیق‌تری نیاز دارد، کاهش یابد. این ارزیابی نشان می‌دهد که ChatGPT با وجود حفاظت‌ها و فیلترهای گسترده، در برابر استفاده مخرب مقاوم نیست: دشمنان می‌توانند از ChatGPT برای ایجاد و استقرار سریع وب‌سایت‌های فیشینگ استفاده کنند، که به طور قابل توجهی خطر بالقوه مرتبط با استفاده از ChatGPT برای چنین فعالیت‌های غیرقانونی را افزایش می‌دهد و گستره و دامنه حملات فیشینگ را گسترش می‌دهد.

سلب مسئولیت از هرگونه سوء استفاده غیر قانونی

با توجه به مطالعات صورت گرفته توسط تیم تحقیقاتی المهدی ما با توجه به امکان مورد سوء استفاده قرار گرفتن محتوای علمی این مقاله با توجه به محرمانگی اطلاعات و پیمان عدم فاش اطلاعات سازمانی ما امکان ارائه اطلاعات ریز پروژه ها را برای تیم پژوهشی خود محفوظ میدانیم و رسماً اعلام میداریم که تمام محتوای این مقاله صرفاً علمی است و از هرگونه استفاده از مطالب این مطالب برای نیت خصمانه صلب مسئولیت می نماییم.

و تمام دامنه های ثبت شده برای بررسی و مطالعه امکان طرح ریزی و پیاده سازی عملی حمله فیشینگ با استفاده از Chat-GPT دامنه های فیشینگ پاکسازی شده و از بین رفتند. فلذا مطالعه ما با مجوز از تارگت مورد حمله بوده است.

مراجع:

- [1] OpenAI, "Introducing ChatGPT," <https://openai.com/blog/chatgpt>.
- [2] OpenAI, "GPT-3.5," <https://platform.openai.com/docs/models/gpt-3-5>.
- [3] P. Ray, "ChatGPT: A Comprehensive Review on Background, Applications, Key Challenges, Bias, Ethics, Limitations and Future Scope," *Internet of Things and Cyber-Physical Systems*, vol. 3, 04 2023.
- [4] Nova, "How ChatGPT is Being Used in Virtual Assistants?" <https://novaapp.ai/blog/chatgpt-virtual-assistant>.
- [5] HelpWise, "The complete Guide to Using Chat-GPT for Customer service," <https://helpwise.io/blog/how-to-use-chatgpt-for-customer-service>.
- [6] M. Daun and J. Brings, "How ChatGPT Will Change Software Engineering Education," in *Proc. of the 2023 Conference on Innovation and Technology in Computer Science Education*, 2023, p. 110–116.
- [7] J. McCoy, "ChatGPT for Content Creation: How to Write More in Less Time," <https://contentatscale.ai/chatgpt-for-content-creation/>.
- [8] Summa Linguae, "How to Use ChatGPT for Translation," <https://summalinguae.com/translation/how-to-use-chatgpt-for-translation>.
- [9] F. Huang, H. Kwak, and J. An, "Is ChatGPT Better than Human Annotators? Potential and Limitations of ChatGPT in Explaining Implicit Hate Speech," in *Companion Proc. of the ACM Web Conference*, 2023.

- [10] CheckPoint, "OPWNAI: Cybercriminals Starting to Use ChatGPT," <https://research.checkpoint.com/2023/opwnai-cybercriminals-starting-to-use-chatgpt/> .
- [11] CyberARC, "Chatting Our Way Into Creating a Polymorphic Malware," <https://www.cyberark.com/resources/threat-research-blog/chatting-our-way-into-creating-a-polymorphic-malware>, 2023.
- [12] Meta, "The Malware Threat Landscape: NodeStealer, DuckTail, and More," <https://engineering.fb.com/2023/05/03/security/malware-nodestealer-ducktail/>, 2023.
- [13] K. T. Gradon, "Electric Sheep on the Pastures of Disinformation and Targeted Phishing Campaigns: The Security Implications of ChatGPT," *IEEE Security & Privacy*, vol. 21, no. 3, pp. 58–61, 2023.
- [14] Y. M. Pa Pa *et al.*, "An Attacker's Dream? Exploring the Capabilities of ChatGPT for Developing Malware," in *Proc. of CSET*. ACM, 2023.
- [15] NameSheild, "ChatGPT, Can You Write a Phishing Email?" <https://blog.nameshield.com/blog/2023/06/15/chatgpt-can-you-write-a-phishing-email/>.
- [16] S. S. Roy, K. V. Naragam, and S. Nilizadeh, "Generating Phishing Attacks using ChatGPT," 2023.
- [17] A. Oest *et al.*, "Sunrise to Sunset: Analyzing the End-to-end Life Cycle and Effectiveness of Phishing Attacks at Scale," in *Proc. of USENIX Security*, 2020, pp. 361–377.
- [18] —, "Inside a Phisher's Mind: Understanding the Anti-phishing Ecosystem Through Phishing Kit Analysis," in *Proc. of eCrime*, 2018.
- [19] P. Zhang *et al.*, "CrawlPhish: Large-scale Analysis of Client-side Cloaking Techniques in Phishing," in *Proc. of IEEE S&P*, 2021.
- [20] S. Maroofi, M. Korczynski, and A. Duda, "Are You Human? Resilience of Phishing Detection to Evasion Techniques Based on Human Verification," in *Proc. of ACM IMC*, 2020, pp. 78–86.
- [21] B. Acharya and P. Vadrevu, "PhishPrint: Evading Phishing Detection Crawlers by Prior Profiling," in *Proc. of USENIX Security*, 2021.
- [22] R. Liu *et al.*, "Inferring Phishing Intention via Webpage Appearance and Dynamics: A Deep Vision Based Approach," in *Proc. of USENIX Security*, 2022.
- [23] S. Maroofi *et al.*, "COMAR: Classification of Compromised versus Maliciously Registered Domains," in *Proc. of IEEE Euro S&P*, 2020.
- [24] A. Blum *et al.*, "Lexical Feature Based Phishing URL Detection Using Online Learning," in *Proc. of ACM AISec*, 2010.
- [25] OpenAI, "Models," <https://platform.openai.com/docs/models/gpt-4>.
- [26] —, "OpenAI Codex," <https://openai.com/blog/openai-codex> .

[27] GitHub, "Your AI Pair Programmer," <https://openai.com/blog/openai-codex> .

[28] A. J. O'Neal, "Chat GPT "DAN" (and other "Jailbreaks"),"

<https://gist.github.com/coolaj86/6f4f7b30129b0251f61fa7baaa881516>.

[29] E. Fasllija *et al.*, "Phish-Hook: Detecting Phishing Certificates Using Certificate Transparency Logs," in *SecureComm*, 2019.

منابع مکمل:

[1]. "How ChatGPT will change cybersecurity", [online] Available:

<https://www.kaspersky.com/blog/chatgpt-cybersecurity/46959/>.

[2]. R. D. Caballar, "Cybercrime Meets ChatGPT: Look Out World", January 2023, [online] Available: <https://spectrum.ieee.org/chatgpt-and-cybercrime>.

[3]. A. Hughes, "ChatGPT: Everything you need to know about OpenAI's GPT-3 tool", January 2023, [online] Available: <https://www.sciencefocus.com/future-technology/gpt-3/>.

[4]. B. Gordijn and H. Have, "ChatGPT: evolution or revolution?", *Med Health Care and Philos*, 2023, [online] Available: <https://doi.org/10.1007/s11019-023-10136-0>.

[5]. F. Salahdine and N. Kaabouch, "Social Engineering Attacks: A Survey", *Future Internet*, vol. 11, no. 4, pp. 89, [online] Available: <https://www.mdpi.com/1999-5903/11/4/89>.

[6]. K. Chetioui, B. Bah, A. Ouali Alami and A. Bahnasse, "Overview of Social Engineering Attacks on Social Networks", *Procedia Computer Science*, vol. 198, pp. 656-661, 2022, [online] Available: <https://doi.org/10.1016/j.procs.2021.12.302>, ISSN 1877-0509.

[7]. L. Irwin, "The 5 Biggest Phishing Scams of All Time", October 2022, [online] Available: <https://www.itgovernance.eu/blog/en/the-5-biggest-phishing-scams-of-all-time>.

[8]. R. Zulfikar, "Phishing attacks and countermeasures" in *Handbook of Information and Communication Security*, Springer, 2010, ISBN 9783642041174.

[9]. A. Kumar Jain and B.B. Gupta, "A survey of phishing attack techniques defence mechanisms and open research challenges", *Journal of Enterprise Information Systems*, vol. 16, pp. 527-565, 2022, [online] Available: <https://doi.org/10.1080/17517575.2021.1896786>.

[10]. B. Gupta, N. Arachchilage and K. Psannis, "Defending against phishing attacks: taxonomy of methods current issues and future directions", *Telecommun Syst*, vol. 67, pp. 247-267, 2018, [online] Available: <https://doi.org/10.1007/s11235-017-0334-z>.

[11] J. Gori Mohamed, M. Mohammed Mohideen, N. Shahira Banu. E-Mail phishing an open threat to everyone. *Int. J. Sci. Res. Publ.* Volume 4, Issue 2, February 2014.

- [12] Postfix. <http://www.postfix.org>.
- [13] Air force mypers. <http://www.afpc.af.mil/Support/myPers/>.
- [14] Google Home Added 600,000 More U.S. Users in 2018 Than Amazon Echo, But Amazon Echo Dot is Still the Most Owned Smart Speaker. <https://voicebot.ai/2019/03/07/google-home-added-600000-more-u-s-users-in-2018-thanamazon-echo-but-amazon-echo-dot-is-still-the-most-owned-smart-speaker/>.
- [15] Alexa top 1 million websites. <https://www.alexa.com/topsites>.
- [16] Alexa top sites by category. <https://www.alexa.com/topsites/category>.
- [17] Alexa Skill Website. <https://goo.gl/PvwL3T>, Amazon Mechanical Turk. <https://www.mturk.com/>.
- [18] Amazon Alexa Skill Count Surpasses 30,000 in the U.S. <https://voicebot.ai/2018/03/22/amazon-alexa-skill-count-surpasses-30000-u-s/>.
- [19] Amazon Alexa Skill Count Surpasses 80,000 worldwide. <https://voicebot.ai/2018/03/22/amazon-alexa-skill-count-surpasses-30000-u-s/>.
- [20] Steps to test Amazon Alexa skill. <https://developer.amazon.com/docs/devconsole/test-your-skill.html#ways-to-test-an-alexa-skill>.
- [21] C. Jones, "The Top 10 Phishing Protection Solutions", January 2023, [online] Available: <https://expertinsights.com/insights/top-10-phishing-protection-solutions/>.
- [22]. V. C. Gungor, D. Sahin, T. Kocak, S. Ergut, C. Buccella, C. Cecati, et al., "A survey on smart grid potential applications and communication requirements", *IEEE Trans. Ind. Informat.*, vol. 9, no. 1, pp. 28-42, Feb. 2013.
- [23]. S. Garg, K. Kaur, G. Kaddoum, J. J. P. C. Rodrigues and M. Guizani, "Secure and lightweight authentication scheme for smart metering infrastructure in smart grid", *IEEE Trans. Ind. Informat.*, vol. 16, no. 5, pp. 3548-3557, May 2020.
- [24]. C. Ang, *The Most Cyber Attsignificant Acks From 2006-2020 by Country*, 2021, [online] Available: <https://www.visualcapitalist.com/cyber-attacks-worldwide-2006-2020/>.
- [25]. M. Benmalek and Y. Challal, "MK-AMI: Efficient multi-group key management scheme for secure communications in AMI systems", *Proc. IEEE Wireless Commun. Netw. Conf.*, pp. 1-6, Apr. 2016.
- [26]. M. Z. Gunduz and R. Das, "Cyber-security on smart grid: Threats and potential solutions", *Comput. Netw.*, vol. 169, Mar. 2020.

- [27]. R. E. Pérez-Guzmán, Y. Salgueiro-Sicilia and M. Rivera, "Communication systems and security issues in smart microgrids", *Proc. IEEE Southern Power Electron. Conf. (SPEC)*, pp. 1-6, Dec. 2017.
- [28]. M. Z. Gunduz and R. Das, "Analysis of cyber-attacks on smart grid applications", *Proc. Int. Conf. Artif. Intell. Data Process. (IDAP)*, pp. 1-5, Sep. 2018.
- [29]. H. Ritchie and M. Roser, *Two-Thirds of Global Population Will Live in Cities By 2050 Our World in Data*, 2018, [online] Available: <https://ourworldindata.org/urbanization>.
- [30]. M. Benmalek, Y. Challal and A. Derhab, "Authentication for smart grid AMI systems: Threat models solutions and challenges", *Proc. IEEE 28th Int. Conf. Enabling Technologies: Infrastructure Collaborative Enterprises (WETICE)*, pp. 208-213, Jun. 2019.
- [31]. A. Halbouni, T. S. Gunawan, M. H. Habaebi, M. Halbouni, M. Kartiwi and R. Ahmad, "Machine learning and deep learning approaches for cybersecuriy: A review", *IEEE Access*, vol. 10, pp. 19572-19585, 2022.
- [32]. A. Hasankhani, S. M. Hakimi, M. Bisheh-Niasar, M. Shafie-khah and H. Asadolahi, "Blockchain technology in the future smart grids: A comprehensive review and frameworks", *Int. J. Electr. Power Energy Syst.*, vol. 129, Jul. 2021.
- [33]. M. Ghiasi, T. Niknam, Z. Wang, M. Mehrandezh, M. Dehghani and N. Ghadimi, "A comprehensive review of cyber-attacks and defense mechanisms for improving security in smart grid energy systems: Past present and future", *Electr. Power Syst. Res.*, vol. 215, Feb. 2023.
- [34]. M. Abdalzaher, M. Fouada, A. Emran, Z. Fadlullah and M. Ibrahim, "A survey on key management and authentication approaches in smart metering systems", *Energies*, vol. 16, no. 5, pp. 2355, Mar. 2023.
- [35]. M. Sewak, S. K. Sahay and H. Rathore, "Deep reinforcement learning in the advanced cybersecurity threat detection and protection", *Inf. Syst. Frontiers*, vol. 25, pp. 589-611, Aug. 2022.
- [36]. S. Banik, S. K. Saha, T. Banik and S. M. M. Hossain, "Anomaly detection techniques in smart grid systems: A review", *Proc. IEEE World AI IoT Congr. (AllIoT)*, pp. 331-337, Jun. 2023.
- [37]. J. Kua, M. B. Hossain, I. Natgunanathan and Y. Xiang, "Privacy preservation in smart meters: Current status challenges and future directions", *Sensors*, vol. 23, no. 7, pp. 3697, Apr. 2023.

- [38]. K. Y. Yap, H. H. Chin and J. J. Klemeš, "Blockchain technology for distributed generation: A review of current development challenges and future prospect", *Renew. Sustain. Energy Rev.*, vol. 175, Apr. 2023.
- [39]. S. S. Koduru, V. S. P. Machina and S. Madichetty, "Cyber attacks in cyber-physical microgrid systems: A comprehensive review", *Energies*, vol. 16, no. 12, pp. 4573, Jun. 2023.
- [40]. M. Lydia, G. E. P. Kumar and A. I. Selvakumar, "Securing the cyber-physical system: A
- [41]. S. Bhattacharjee and S. K. Das, "Detection and forensics against stealthy data falsification in smart metering infrastructure", *IEEE Trans. Depend. Secure Comput.*, vol. 18, no. 1, pp. 356-371, Jan. 2021.
- [42]. X. Lou, "Learning-based time delay attack characterization for cyber-physical systems", *Proc. IEEE Int. Conf. Commun. Control Comput. Technol. Smart Grids (SmartGridComm)*, pp. 1-6, Oct. 2019.
- [43]. A. Ameli, A. Ayad, E. F. El-Saadany, M. M. A. Salama and A. Youssef, "A learning-based framework for detecting cyber-attacks against line current differential relays", *IEEE Trans. Power Del.*, vol. 36, no. 4, pp. 2274-2286, Aug. 2021.
- [44]. S. Yankson and M. Ghamkhari, "Transactive energy to thwart load altering attacks on power distribution systems", *Future Internet*, vol. 12, no. 1, pp. 4, Dec. 2019.

The other side of the coin: the dark side of artificial intelligence; Practical cyber attack with artificial intelligence

Sina Manafi¹, Fardin Esmaili Sangri², Mehrdad Akbari³

Security Operations Center Researcher,

Sinahgfv@gmail.com

Computer Department, Urmia Branch, Islamic Azad University, Urmia, Iran,

Fardin.e.s62@gmail.com

Bachelor of Information Technology Engineering - Information Technology -

Comprehensive University of Applied Sciences - Khaneh Kargar University

Education Center, Tabriz, Tabriz, Iran, Email

Abstract— This paper explores the possibility of using ChatGPT to develop advanced phishing attacks and automate large-scale deployments. We make ChatGPT implement the following parts of a phishing attack: 1) impersonate a target website, 2) integrate code to steal user information, 3) obfuscate code, 4) automatically deploy the website to your hosting provider, 5) register phishing domain, and 6) integrating the website with a reverse proxy. It will make the initial evaluation of automatically generated phishing kits practical. It highlights their rapid production and deployment process as well as the close similarity of the resulting pages to the intended website. More broadly, we show that recent advances in artificial intelligence confirm the potential dangers of its misuse in phishing attacks, which could lead to an increase in their prevalence and severity. This highlights the need for advanced countermeasures and defensive measures in artificial intelligence systems.

Keywords: cyber attacks, artificial intelligence, ChatGPT, automation of attacks