

Desafíos sobre las nuevas tecnologías de resolución de CAPTCHA y características de evolución de CAPTCHA en el futuro próximo

Daniel Alejandro Maldonado Ruiz (IEEE), Juan Alejandro Devincenzi (asesor)

Universidad de Buenos Aires, Facultades de Ingeniería, Ciencias Exactas y Naturales y Ciencias Económicas, Escuela de Estudios de Postgrado. Av. Córdoba 2122 (C1120AAQ), Ciudad Autónoma de Buenos Aires. Argentina.

Autores para correspondencia: daniel.a.maldonado@ieee.org, juanalejandro.devincenzi@gmail.com

Fecha de recepción: 28 de septiembre 2015 - Fecha de aceptación: 12 de octubre 2015

RESUMEN

Desde el desarrollo de los sistemas de computación, se consideró necesario diferenciar las respuestas que provenían tanto de un sistema informático como de una persona, para que uno no pueda suplantar al otro en ambientes donde sólo una persona deba utilizar, lo cual se realizó a través de CAPTCHA, que es un eficiente arreglo desde el punto de vista criptográfico, pero que al paso del tiempo se ha encontrado con los nuevos sistemas de Inteligencia Artificial. Se incluye un breve análisis teórico de las razones de su uso desde el punto de vista de la criptografía, de la evolución del sistema en su presentación para los usuarios, haciendo énfasis en los CAPTCHA más fuertes del mercado. Adicionalmente una revisión de la tecnología que ha permitido vulnerar a cierto nivel al sistema y la evolución de los sistemas de Inteligencia Artificial que a mediano plazo podrían vulnerar CAPTCHA; para finalmente realizar recomendaciones de cómo se debe diseñar y mantener CAPTCHA fortalecido al corto y mediano plazo frente a la evolución de la tecnología y que características debería tener para que pueda seguir siendo un test válido para diferenciar un humano de una máquina.

Palabras clave: CAPTCHA, inteligencia artificial, criptografía, algoritmos de reconocimiento óptico, seguridad, usabilidad, Internet.

ABSTRACT

Since the development of computer systems, it becomes necessary to differentiate the responses came from both a computer system as a person, so that one cannot replace the other in environments where only one person should use, which was made through CAPTCHA, which is an efficient arrangement from the cryptographic viewpoint, but it has found several issues with the new systems of Artificial Intelligence. This Paper includes a brief theoretical analysis of the reasons for their use from the cryptography viewpoint, the evolution of the system in its presentation for the users, focusing on the strongest CAPTCHA in the market. Additionally, a review of the technology that has allowed to break the system at certain level and the evolution of Artificial Intelligence systems in the medium term could defeat CAPTCHA; to finally make recommendations on how to design and maintain a strong CAPTCHA at short and medium term against the evolution of technology and what kind of features should have CAPTCHA to remain a valid test to differentiate a human from a machine.

Keywords: CAPTCHA, artificial intelligence, cryptography, optic recognition algorithms, security, usability, Internet.

1. INTRODUCCIÓN

Desde el desarrollo de los sistemas de computación, se consideró necesario diferenciar las respuestas que provenían tanto de un sistema informático como de una persona, para que uno no pueda suplantar

al otro. Alan Turing diseñó un test que permitía realizar esta comprobación, mediante preguntas que solo los humanos estábamos en capacidad de responder adecuadamente, es decir servía para identificar máquinas que puedan pensar como humanos.

Ya en los primeros días de Internet, la necesidad de crear textos que sean ilegibles a las computadoras formó parte de los estudios de los desarrolladores. En 1997 Altavista diseñó un sistema para evitar que bots ingresaran url's a su motor de búsquedas. Pero no sería hasta 2000 cuando investigadores de la Universidad Carnegie-Mellon e IBM desarrollaron un test para evitar que bots desde Internet pudieran acceder y monopolizar servicios para usuarios. Este test se denominó CAPTCHA, o Prueba de Turing Automática y pública para diferenciar humanos de máquinas, por sus siglas en inglés (Completely Automated Public Turing test to tell Computers and Humans Apart).

El presente trabajo pretende dar una visión de cuál es la fortaleza de CAPTCHA pensado desde su origen criptográfico, de cómo los avances en reconocimiento gráfico han permitido vulnerar de alguna manera este sistema y de que características debería poseer CAPTCHA en el futuro para poder seguir siendo seguro frente a cualquier ataque, sea este algorítmico o de inteligencia artificial; es decir, que aun solo pueda ser resuelto por humanos.

En la Sección 2, se explica brevemente cómo un grupo de conceptos criptográficos pueden ser inferidos para la creación de CAPTCHA, en la Sección 3 cuales han sido las vulnerabilidades encontradas mientras el sistema ha sido implementado. En la Sección 4 se especifican los avances en diferentes tecnologías de evolución de Inteligencia Artificial, y en la Sección 5 las recomendaciones presentes y futuras para la aplicación de CAPTCHA a mediano y largo plazo. Finalmente la Sección 6 contiene las conclusiones fruto de este trabajo.

Es importante comprender que aun cuando un sistema pueda ser visto como infalible al corto plazo; la tecnología de Inteligencia Artificial crece cada vez más y es un desafío mantener los sistemas seguros en un mundo donde las máquinas están tomando control de comunicaciones y servicios con los que nos relacionamos día a día.

2. CAPTCHA COMO MODELO DE SEGURIDAD

CAPTCHA fue diseñado para generar y verificar test que una persona podría fácilmente resolver, pero que a un sistema computacional le sea muy difícil. De acuerdo al diseño original, ni siquiera el generador del test en sí mismo está capacitado para pasarlo. Esta particularidad del diseño ha hecho que efectivamente sirva para crear una brecha entre humanos e inteligencia programable (von Ahn, 2003). Esta capacidad de diferenciación ha posibilitado que CAPTCHA sea utilizado como:

- Verificador de votantes en encuestas.
- Comprobador de identidades para la creación de correos electrónicos y demás servicios sociales en la Red.
- Implementar restricciones de acceso frente a motores de búsqueda.
- Defensas frente al envío masivo de correo basura y malware a través de listas automáticas.
- Defensa frente a ataques de diccionario en contraseñas de sistemas informáticos (Pinkas, 2002).

El concepto público de CAPTCHA implica que el algoritmo de creación e implementación debe ser público, con excepción de la porción de código que maneja la aleatoriedad para su utilización. Este concepto hace que mientras las vulnerabilidades de otros sistemas de seguridad están basadas en la ruptura de un algoritmo específico, la ruptura de CAPTCHA es un problema específico de Inteligencia Artificial, en donde encontrar un sistema que pueda batir a CAPTCHA eficientemente es encontrar un sistema de Inteligencia Artificial altamente eficiente y confiable.

Hablando desde un punto de vista sistemático, no se puede afirmar que CAPTCHA sea imbatible programáticamente dado que un programa, como el funcionamiento del cerebro humano, puede resolverlo. Sin embargo, los creadores presentaron evidencia matemática en la cual se afirma que escribir un programa que pueda romper CAPTCHA es altamente complicado, afirmando que: “Un

CAPTCHA es un protocolo criptográfico cuya dificultad asumida está basada en un problema de Inteligencia Artificial” (von Ahn, 2003).

Moni Naor, del Instituto Weismann de ciencias de Israel, en 1997 desarrolló la idea de los test automáticos de Turing (Naor, 1997), que dio pie a los primeros desarrollos del sistema presentados por Altavista; y en general este sistema fue útil hasta el desarrollo de los sistemas de reconocimiento óptico de caracteres (Optical Character Recognition u OCR), fruto de lo cual se desarrolló CAPTCHA como un Test de Turing que no esté basado en la dificultad presentada por los OCR's.

3. FACTORES DE VULNERABILIDAD DE CAPTCHA

El teorema establecido en la publicación original de CAPTCHA (von Ahn, 2003), si bien define los parámetros de un CAPTCHA práctico, sólo define el cómo debe ser escrito el programa, no con qué características debe ser presentado al usuario para que sea más o menos vulnerable a los nuevos sistemas de OCR.

En años recientes, y en función a estudios realizados sobre el sistema de seguridad de CAPTCHA¹ (Bursztein, 2010) determinaron que con un adecuado sistema de segmentación de la imagen a ser atacada, es posible encontrar y vulnerar determinados tipos de sistemas según su tipo de construcción, caracteres y distorsión; sin embargo, la segmentación por sí sola no es suficiente en un entorno variable de imágenes en internet.

Estos estudios (Bursztein, 2011) determinaron que para poder vulnerar CAPTCHA y romperlo según su complejidad se requieren de 5 pasos genéricos:

- Pre-Procesamiento
- Segmentación
- Post-Segmentación
- Reconocimiento
- Post-Procesamiento

Por segmentación estaba entendida la separación de una secuencia de caracteres en caracteres individuales, y por reconocimiento la identificación de los mismos, pero los nuevos sistemas han incorporado arreglos adicionales a los caracteres distorsionados como tal, que interfieren con los pasos básicos, por lo que los nuevos pasos previos aseguran la mayor normalización de los caracteres para su posterior reconocimiento. Finalmente el post-procesamiento puede mejorar la precisión de los caracteres para facilitar, por ejemplo, las funciones de aprendizaje de los sistemas atacantes, o determinar si los conjuntos de caracteres son formados al azar o si corresponden a palabras reales en un idioma determinado.

Este procedimiento se derivó de un estudio en el que se comprobaba la dificultad que tenían las personas para resolver CAPTCHAs que se utilizan en los sitios más visitados según Alexa², además de sistemas proporcionados por recaptcha.net y captchas.net. Este estudio determinó que los humanos, en promedio, solo acertaban a resolver CAPTCHA un 71% de las veces³ (Bursztein, 2011), que las personas que no tienen un nivel fluido de inglés (idioma en el que se presentan la mayoría de CAPTCHA) tienen mayor problema que los que son angloparlantes, y que, con la edad, las personas se vuelven más precisas en resolver un CAPTCHA, aunque sacrificando la velocidad de resolución. Mientras más respuestas diferentes posea un CAPTCHA para los testadores, mayor es la posibilidad de fallo, y por tanto, mayor dificultad presenta para las personas.

Sin embargo, se demostró que lo que era difícil de leer para una persona no era necesariamente difícil de interpretar para un sistema de reconocimiento visual y viceversa. Los algoritmos planteados

¹ En este trabajo no se consideran los ataques manuales a CAPTCHA, por parte de personas contratadas para este propósito.

² <http://www.alex.com/topsites> (visto el 08/09/15).

³ A pesar de que aquel estudio consideró CAPTCHAs de audio, este trabajo está enfocado en CAPTCHA de lectura, por lo que los resultados obtenidos no serán tomados en cuenta.

para resolver CAPTCHA se hicieron de acuerdo a las técnicas que usan las personas, aunque considerando aquellos que tenían mayor número de respuestas diferentes.

Para evaluar la efectividad real de los ataques propuestos, el primer paso es medir la precisión, es decir, la fracción de CAPTCHA que pueden ser resueltas adecuadamente. Cada sistema, sin embargo, según su constitución y diseño presentará diferentes resultados a los análisis y adicionalmente, dependerá de las restricciones propias del sitio en cuestión y de la confianza del atacante para resolver determinado CAPTCHA sin cambiarlo por otro más fácil.

La cobertura se define como la cantidad de CAPTCHA que el atacante está dispuesto a resolver, y la precisión como la cantidad de CAPTCHA respondidos correctamente. En los apartados anteriores, y a partir del teorema señalado, se estableció como CAPTCHA exitoso a aquel que no puede ser vulnerado más de 1 de 10000 intentos, es decir en un porcentaje del 0.01%. En la realidad se puede considerar un CAPTCHA vulnerado si la precisión puede alcanzar al menos el 1%.

El proceso definido para vencer un CAPTCHA según experimentación previa es el de asignar un tipo de segmentación específica para cada sistema con un sistema de aprendizaje basado en un OCR (Bursztein, 2011). Primero, el programa pre-procesa la imagen para que pueda ser fácil de analizar, removiendo colores solidos de fondo y aplicando filtros reductores de ruido. Luego, el programa trata de separar la imagen en trozos que contengan un carácter, o un porcentaje muy elevado de uno, finalmente, estos trozos son vectorizados para análisis o ingresados en una red neural para determinar qué carácter está almacenado en cada uno de los trozos y dar una respuesta del CAPTCHA analizado. Los experimentos mencionados han demostrado que mientras más fácil sea segmentar un CAPTCHA más vulnerable es. Los diferentes CAPCHAs en el mercado han establecido sus propias técnicas anti-segmentación para evitar que se puedan separar en caracteres y analizar el contenido, con mayor o menor éxito, ya que solo son efectivas si se han diseñado e implementado adecuadamente, además que deben ser implementados en todas las capas de la imagen, a fin de crear un esquema seguro.

El reconocimiento de caracteres es un problema de fondo dentro del análisis para encontrar una forma eficiente de romper CAPTCHA. La base de datos de caracteres manuscritos del Instituto Nacional de Estándares y Tecnología de Estados Unidos (Mixed National Institute of Standards and Technology, o MNIST, por sus siglas en inglés) es la principal fuente de entrenamiento para sistemas basados en reconocimiento de caracteres distorsionados y de cómo reconocerlos eficientemente.

Cuando CAPTCHA no puede ser segmentado y se debe hacer un reconocimiento sobre toda una imagen, es posible realizar aproximaciones alternativas con descriptores de imagen complejas, como SURF y SIFT⁴ (Bursztein, 2011), que funcionan invariablemente ante la rotación de caracteres y son muy estables contra la distorsión, gracias a su uso de ‘puntos de interés’ que permiten una aproximación mucho más estable y rápida. Como el número de ‘puntos de interés’ no puede ser normalizado porque varía según cada tipografía, no se pueden aplicar clasificadores para hacer más eficiente el reconocimiento como tal.

Para aumentar la seguridad de anti-segmentación, se utiliza también fondos para las imágenes a ser descifradas que contribuyan a la confusión de la imagen general al mezclar el texto con las características del fondo. La idea de usar fondos complejos, parte de que, usando sus líneas/formas se puede contribuir a la confusión de un detector de caracteres que impida la segmentación de los mismos. A pesar de que estudios previos de detección demostraron que usualmente este tipo de defensa es insegura, muchos CAPTCHA en internet la siguen usando, valiéndose de fondos aleatorios generados en función de patrones pre-establecidos que buscan evitar que los atacantes consigan el patrón de uso de las imágenes, mientras que mantienen destacados los caracteres para que los usuarios puedan resolverlos.

El problema radica en que las imágenes utilizadas como fondos no deben ser de alta calidad para que no comprometan el peso del programa como tal, por lo que necesariamente y aun a pesar de su aparente complejidad, el número de colores utilizados en cada una se mantiene finito; sin mencionar el color con el que se destaca cada carácter. A partir de estas características los analizadores remueven todos los pixeles que sean diferentes a los del color de los caracteres, si se ha comprobado que estos

⁴ SURF: Características robustas de alta velocidad o Speeded-Up Robust Features.

SIFT: Transformada de Características de Escala Invariante o Scale Invariant Feature Transform.

tienen el suficiente número de píxeles para ser considerado como tal. Otro uso de los fondos es el de asignarles un color muy similar al de los caracteres para que estos puedan perderse en el fondo. El problema está en que muchos colores que para las personas no son tan cercanos, si lo son a nivel del espectro RGB⁵; y viceversa. Sin embargo, una forma efectiva de contrarrestar consiste en romper el patrón con una diferente representación de colores que este más cercana a la percepción humana, y luego binarizar el CAPTCHA basado en el matiz o en la saturación del mismo para facilitar el trabajo de los analizadores (Bursztein, 2010).

Añadir ruido aleatorio como técnica de confusión y anti-segmentación está considerada como la manera más eficiente de asegurar un CAPTCHA, pero el mismo tiene que ser del mismo color de los caracteres o fácilmente los reconocedores de patrones pueden eliminarlo. Para eliminar el ruido de un CAPTCHA, se han utilizado varias técnicas de filtrado con los años, pero la más eficiente es usar un Campo Aleatorio de Markov, también conocido como algoritmo Gibbs.

Gibbs es un algoritmo iterativo que computa la carga de color de cada píxel, o su energía, basándose en los píxeles que tiene alrededor para remover los píxeles cuya energía este por debajo de un umbral predeterminado, terminando cuando no hay más píxeles por remover. La energía de cada píxel es calculada resumiendo los valores en una escala de grises de los 8 caracteres circundantes y dividiéndola para 8. Este tipo de algoritmo hace que la mayoría de fondos utilizados, con o sin ruido, sean inseguros en un ambiente de Internet.

Otra técnica para evitar la segmentación es el uso de líneas que crucen los caracteres que conforman el CAPTCHA, ya sea con líneas pequeñas a través de determinados caracteres o líneas largas que atraviesan todo el CAPTCHA. Para cuando se utilizan líneas cortas entre los caracteres, el estándar es utilizar una segmentación basada en histogramas que proyecta a los píxeles de CAPTCHA dentro de coordenadas X o Y. Esta aproximación funciona porque las regiones donde hay mayor densidad de caracteres se crean picos en el histograma. El problema consiste en la determinación del umbral y el tamaño de las ventanas alrededor del mismo; para evitar las líneas pequeñas es mucho más efectivo utilizar el algoritmo de Gibbs con reconstrucción de caracteres y no requiere de una afinación tan intensiva como los histogramas.

Para cuando se utilizan líneas largas es común utilizar líneas que tengan el mismo grosor de los segmentos de los caracteres, lo que hace que no sean vulnerables a las técnicas anti-ruido, como Gibbs, pero es susceptible frente a filtros buscadores de líneas, como las transformadas de Hough o la detección de borde de Canny, que encuentran líneas sobre un texto de manera muy eficiente. La dificultad radica en no deformar los caracteres cuando las líneas son removidas, lo que se logra comprobando los píxeles circundantes a uno con posibilidad de remoción para determinar si en efecto debe ser removido o no. Esta técnica se dificulta cuando los caracteres son huecos, lo que hace que queden deformados más allá de toda recuperación cuando las líneas que han sido superpuestas son removidas.

Finalmente (Bursztein, 2011), la unión de caracteres o colapso es considerada como la mejor y más segura técnica de anti-segmentación, dependiendo solamente de sí su uso al generar las imágenes es correctamente aplicado. Esta es la razón por la que se distinguen dos tipos de colapso: uno donde el atacante puede predecir la segmentación de caracteres a pesar de encontrarse colapsados y otro donde no se puede predecir y el atacante debe realizar ataques de ‘fuerza bruta’ al sistema.

El colapso predictivo se da cuando a pesar de los caracteres están juntos, el atacante aún puede definir donde están las separaciones de cada uno para la segmentación, sobre todo si los caracteres son lo suficientemente regulares en su construcción o si el número de los mismos viene predefinido. Conocido como segmentación oportunista, ya que se basa en información indirecta para realizar su trabajo. Es la más vulnerable de las uniones de caracteres.

El colapso no predictivo se da cuando el número de caracteres no es conocido por defecto y cada carácter tiene su propio tamaño y tipografía, lo que impide obtener patrones definidos para segmentar y obliga a los atacantes a tratar de reconocer el CAPTCHA en su totalidad sin segmentación. Ante

⁵ RGB (Red, Green, Blue): modo de representación de imágenes en la cual cada píxel es una combinación de los tres colores mencionados, cada uno representado en 8 bits.

este tipo de CAPTCHA es que se han utilizado plantillas de caracteres manuscritos a partir de la MNIST o el uso de redes neurales para el reconocimiento.

4. AVANCES EN INTELIGENCIA ARTIFICIAL

4.1. *Redes neuronales convolucionales*

Las operaciones que para el cerebro humano son triviales no lo son para las máquinas, ya que el funcionamiento secuencial y paralelo del procesamiento del cerebro requiere obligatoriamente un elevado número de interconexiones masivas y nodos de interpretación, en este caso neuronas y sinapsis.

Los sistemas que más se han acercado a emular el sistema de procesamiento del cerebro son las redes convolucionales con sistemas y algoritmos de entrenamiento basados en gradiente. Sin embargo, la desventaja de estas redes está en que su emulación de los procesos cerebrales se limita a reconocer fotogramas, es decir, en imágenes estáticas escaneadas y no están diseñados para procesar eventos en tiempo real. De hecho, se ha comprobado que estos sistemas son más eficientes en reconocimiento que en segmentación, ya que el número de conexiones necesarias aun es una limitante física al momento de su construcción.

En este sentido, en 2005 un grupo de investigadores de Microsoft determinó (Who Made That Captcha?, 2014) que cuando se ha resuelto el problema es la segmentación, se puede asumir directamente que también se ha resuelto el problema del reconocimiento. Para este objetivo, se propusieron comprobar la eficiencia de una red convolucional en el reconocimiento de caracteres singulares y evaluar su rendimiento frente a la capacidad de reconocimiento de un ser humano. Por ello diseñaron un experimento en el cual generaron caracteres aleatorios y los deformaron según secuencias computadas para que puedan ser consideradas como parte de la base de datos de la MNIST y por intervalos se evaluaron según su complejidad por la red neuronal y por los humanos, y comprobaron que, en este escenario controlado, las redes neuronales son mucho más eficientes en el reconocimiento que los humanos, comprobando además que el reconocimiento no es un problema una vez que se ha salvado el paso de la segmentación.

Este experimento aún se basaba en imágenes estáticas para realizar las comprobaciones con las redes neuronales, hasta que en 2011, un grupo de investigadores de la Universidad de Sevilla (Goodfellow, 2013) diseñó una red convolucional bio-inspirada, a través de chips convolucionales programables que manejan el protocolo AER (Representación de Eventos mediante Direcciones, o Address-event representation) para el reconocimiento ya no de fotogramas, sino de eventos en tiempo real.

Este sistema de 6 capas emula el funcionamiento de las neuronas donde cada neurona en un mapa específico está conectada únicamente con neuronas de la capa siguiente. Esta característica hace que para cada entrada de un evento se separan por etapas, cada una con diferentes imágenes de salida conocidas como ‘mapas de características’, que están compuestas de los mencionados mapas de neuronas.

Para su experimento realizaron una simulación de sistemas AER donde multiplexaron cada evento de entrada en 6 canales conectados a la primera capa de módulos AER, que implementan filtros de Gabor para el análisis. A la salida de la primera capa se envían a la siguiente donde se submuestran y se re-codifican las direcciones de las entradas para obtener unas que las siguientes capas puedan reconocer. Cada una de estos canales almacena la dirección específica de una parte del evento y se replican en 4 canales adicionales, que envían la información a la tercera capa que no es más que una nueva estructura de convolución con 6 puertos de entrada, donde la máscara de convolución que corresponde a cada entrada se añade a la dirección previamente codificada en la matriz de píxeles.

En función de umbrales y tiempos definidos, los eventos de salida son enviados a la capa siguiente para su procesamiento mientras que la previa es reseteada. Estos tiempos refractarios son utilizados para emular las no-linealidades de los modelos habituales de redes neuronales. En la cuarta

etapa se vuelve a submuestrear los eventos resultantes y a enviarse a la quinta etapa donde son replicados 10 veces y enviados a una sexta capa compuesta por analizadores i&f (integración y disparo), que presentan resultados positivos o negativos en función de la correspondencia entre la entrada y su salida. El uso de las no-linealidades hace que se presenten solo valores positivos en el que la salida corresponda al dígito de entrada.

A través de esta configuración de los chips AER y de las no-linealidades lograron una tasa de reconocimiento del 93%, lo que minimiza los tiempos de respuesta al obtener un primer resultado a la salida del analizador, que en este caso fue de 4.3 [μ s]. También comprobaron que cuando el flujo de entrada variaba entre dígitos diferentes el tiempo para una respuesta positiva luego de la transición es de 22.4 [μ s]. Sin embargo, y a pesar que la evolución del hardware hacia nuevos sistemas CMOS/noCMOS, aún no ha logrado resolver el problema de la correcta segmentación de caracteres para su reconocimiento.

Adicionalmente, y fruto de StreetView, Google se hizo con miles de números de direcciones de todas las calles a las cuales su sistema Maps ha llegado. A través de ReCaptcha, los usuarios han podido descifrar muchos de los números de direcciones, y, siguiendo este patrón, ha diseñado una red neural para leer los números adquiridos por Street View sin necesidad del componente humano (Google's Street View, 2014).

Esta red realiza un reconocimiento completo de la imagen, emulando el funcionamiento del cerebro humano en lugar de dividirlo en números iguales para su procesamiento. Esta funcionalidad, entrenada a través del Conjunto de Datos de Números de Casas de Street View (SVHN, o Street View House Numbers, por sus siglas en inglés) en algo más de 6 días tiene un nivel de precisión de al menos el 96% frente a la precisión conseguida por los seres humanos a través de ReCaptcha, que es del 98%, el cual los ingenieros de Google se han puesto como el umbral de éxito.

Para su cometido emplearon una red convolucional avanzada, basada en la implementación de DistBelief (2012) que permite entrenar redes neuronales avanzadas distribuidas para analizar imágenes de alta calidad. A través de arquitectura muy avanzada han logrado mejorar el rendimiento con 11 capas ocultas.

La particularidad del uso de esta red está en que las fotografías tomadas por Street View, aunque no presentan letras colapsadas, según el ángulo de fotografía, la luz, el tiempo meteorológico y la forma de la nomenclatura de cada casa hacen que su análisis no sea tan simple, sobre todo si se lo va a realizar sin segmentación.

Para conseguir el nivel mencionado de precisión, realizaron experimentos preliminares con caracteres únicos desplazados aleatoriamente para entrenar a la red en el reconocimiento a través de aumentar la base de datos de estudio en un 30%, aumentando las variables, como la escala de las imágenes, haciendo que un dígito pueda entrar en un contenedor definido, mientras que más caracteres debieron ser reducidos para que puedan caber en los contenedores definidos. Sin realizar este previo entrenamiento, comprobaron que se perdía casi la mitad de precisión en puntos porcentuales.

La arquitectura del sistema consiste en 8 capas convolucionales ocultas, una adicional conectada localmente y dos conectadas densamente entre sí, todas conectadas de tal modo que la salida de la una invariablemente iba a la siguiente sin perder conexiones intermedias. Mientras que la primera contiene las unidades de estudio completas las siguientes están constituidas de rectificadores, dando lugar a 3702 unidades entre todas las capas conectadas todas ellas normalizadas y configuradas para preservar su tamaño de representación.

Google aspira con estos avances en inteligencia artificial poder leer y almacenar direcciones de una manera más eficiente, y extrapolar los resultados a señalética urbana y anuncios publicitarios. Todos estos avances, sin embargo, siguen sin resolver el problema de la segmentación.

4.2. Redes neuronales recursivas corticales

El desarrollo de las redes corticales está íntimamente relacionado con el modo en que neurológicamente los receptores y transmisores envían las señales que los ojos reciben para que sean interpretados en los centros cerebrales de reconocimiento.

Este es un paso intermedio entre la neurología y la Inteligencia Artificial, ya que intenta simular el método en que la experiencia visual consiente a través de áreas corticales interconectadas por medio de un flujo combinado de estímulos y el resultado integral de las respuestas a estos estímulos y de toda la actividad subyacente en estas áreas corticales de la corteza occipital, en donde se encuentran almacenados los receptores de visuales, es decir, una red neuronal recursiva (Chellapilla *et al.*, 2005).

El principio de una Red Neuronal Cortical es el desarrollo de un framework sistemático para determinar el funcionamiento de los bucles recurrentes de las unidades corticales que funcionan como foto receptores y foto transmisores y que son básicas para producir la experiencia de fenómenos visuales consientes asimilados por el sistema nervioso (Pérez-Carrasco, 2011).

Parte del estudio neurológico de estas redes y la parte más difícil de emular por parte de los sistemas computacionales son las representaciones neuronales del sentido de si-mismo como filtro de las representaciones sensoriales que las redes corticales necesitan como inicio de la experiencia consiente y de cómo los bucles de información sobrealimentada en los niveles corticales añaden experiencias sensoriales. Sin embargo para su aplicación solo se toman aproximaciones de su funcionamiento el cual pueda ser emulado mediante algoritmos.

Las redes corticales están basadas en un modelo de aprendizaje de maquina conocido como Memoria temporal jerárquica (Hierarchical Temporary Memory o HTM, por sus siglas en inglés) (Chellapilla *et al.*, 2005) que recoge estas aproximaciones en el funcionamiento de los sistemas occipitales y las propiedades algorítmicas resultantes del neo córtex como un método de presentar un modelo mucho más complejo de la actividad sensorial cerebral.

La idea de su implementación se basa es someter a la red a variados flujos sensoriales en vez de a estímulos propios de respuestas programadas, dando lugar a una red matricial claramente dividida en sub-espacios donde se almacena y procesa la información recibida en función de la 'memoria' que contenga la red, es decir, el entrenamiento al que ha sido expuesta previamente.

La jerarquía de estas redes está determinada en función del tiempo y se establece en tres niveles de jerarquía compuestos por nodos específicos. Mientras más elevada la jerarquía, requiere de menos nodos, ya que reutilizan información adquirida en niveles previos lo que contribuye al mejor procesamiento de patrones complejos pero que limita el procesamiento espacial en niveles superiores. Cada uno de estos nodos tiene básicamente la misma funcionalidad, y se diferencian en la forma del procesamiento de la información que reciben y de las secuencias temporales en las que se producen.

El método de aprendizaje distribuido hace que los estímulos solo interactúen con determinados nodos activos en momentos dados, y aumentando o reduciendo el volumen de nodos involucrados en tanto el nivel de complejidad del estímulo de entrada. Cada uno de estos bloques de nodos activados distribuidamente están interconectados fuertemente entre si simulando las capas del neo córtex, donde las células son capaces de recordar determinados estados previos, pudiendo estar en modo activo, pasivo o predictivo.

Estos tres modos de cada nodo le permiten a la red poder emular predicciones en función de las conexiones que han sido creadas por las entradas previas, enseñado a los bloques sobre cuales deben permanecer encendidos o apagados según evoluciona el flujo de entrada, lo que hace que a la salida de cada sección existan nodos que estén al mismo tiempo en modo activo y predictivo, proporcionando estabilidad temporal cuando el flujo de entrada consiste en patrones largos.

Estos algoritmos de aprendizaje son capaces de mantener un nivel de admisión de datos continuo mientras continúa el flujo de datos, lo que hace que el algoritmo necesite inferir si los datos recibidos están en secuencias previamente aprendidas para evitar el doble procesamiento y avisar a los nodos que contienen esa información que la adicione a la salida del nivel, reduciendo los costos de procesamiento y aumentando la velocidad de aprendizaje en niveles superiores, además de llenar patrones faltantes e interpretar datos ambiguos o de difícil interpretación.

Bajo estos conceptos, la startup de Inteligencia Artificial Vicarius FPT Inc. (AI Startup Vicarious Claims Milestone In Quest To Build A Brain: Cracking CAPTCHA, 2013) (Captcha test 'cracked' by US firm Vicarious, 2013) (Software Firm Claims Breakthrough in Computer Vision Will Lead to Better AI, 2013) utilizó los conceptos de las redes corticales para diseñar una aproximación matemática del funcionamiento del cerebro humano a partir de lo que esta empresa consideraba fallas

de procesamiento de información en los niveles superiores, que impedían un aprendizaje escalado eficiente y un correcto tratamiento de datos, sobre todo en lo que se refería a multimedios.

A través de estos algoritmos, la empresa californiana afirma que emula el funcionamiento de reconocimientos de patrones del cerebro humano, haciendo que la segmentación y el reconocimiento ocurran en el mismo proceso sin realizar una limpieza previa del CAPTCHA en cuestión, resolviéndolo en poco tiempo y eliminando la barrera que hacía a CAPTCHA seguro, tal como lo muestran en su video demostrativo⁶, donde afirman comprobar que su sistema algorítmico puede romper con un éxito del 90% todos los CAPTCHA del mercado, incluyendo los más seguros de Google y ReCaptcha (Vicarious AI passes first Turing Test: CAPTCHA, 2013) (Artificial Intelligence Breaks CAPTCHA Protection, Coders Claim, 2014) (Vicarious AI breaks CAPTCHA ‘Turing test’, 2014).

Sin embargo, y a pesar de estos resultados aparentemente sólidos, Vicarius se ha negado a publicar su investigación a través de una publicación académica acreditada, por lo que sus resultados no han podido ser evaluados por entidades independientes. Luis von Ahn, de la Universidad Carnegie-Mellon y co-creador de CAPTCHA no siente que el sistema, al menos al corto plazo se vea amenazado. En una entrevista con relación a los resultados de Vicarius, afirmó que “es la décima o doceava vez que alguien clama haber roto CAPTCHA”. Estas circunstancias han hecho, que a pesar de la tecnología supuestamente desplegada, una herramienta que pueda derrotar a CAPTCHA no sea viable en el futuro próximo.

5. EVOLUCION Y MANUTENCION DE CAPTCHA

5.1. Características para mantener CAPTCHA seguro al corto plazo [5]

Los estudios sobre seguridad de CAPTCHA determinaron las características que estos sistemas comparten, tanto en creación como en vulnerabilidades. A partir de estos estudios, y en función de todas las características de anti-segmentación y anti-reconocimientos, se han condensado ciertas características y recomendaciones que debe seguir CAPTCHA para poder mantenerse seguro en internet y mantener las características del teorema expuesto previamente.

La construcción de un CAPTCHA seguro parte desde el diseño del núcleo y de las características de anti-segmentación y anti-reconocimiento. La anti-segmentación solo funcionará adecuadamente si el núcleo y el anti-reconocimiento están correctamente aplicados. Si uno de los niveles falla, el atacante ampliará sus probabilidades de vencer al sistema.

Para estas tres capas se han propuesto principios que en su conjunto permiten la creación de CAPTCHAS de texto más seguros:

- Principios de diseño de núcleo:
 - Aleatorizar la longitud de la palabra utilizada como CAPTCHA, para evitar que el atacante pueda usar el número de caracteres como herramienta.
 - Aleatorizar el tamaño de cada carácter, para evitar que los algoritmos de reconocimiento puedan utilizar los mismos puntos de reconocimiento normalizado en cada carácter, haciendo más confusa la lectura para el lector.
 - ‘Ondular’ los CAPTCHA, para aumentar la dificultad de encontrar los puntos de corte en caracteres colapsados o uso de líneas.
 - Evitar usar fondos, ya que los pre-procesadores pueden eliminarlos sin problema y no aportan a la seguridad del sistema.
- Principios de anti-reconocimiento:
 - Fortalecer los esquemas de anti-reconocimiento aplicados en el núcleo, para reducir la posibilidad de detección ante analizadores cada vez más sensibles.

⁶ <https://vimeo.com/77431982>

- No usar un complejo esquema de caracteres, ya que los caracteres especiales y otros símbolos no aumentan la seguridad del esquema, sino solo hacen que sea más difícil para las personas resolverlos.
- Principios de anti-segmentación:
- Usar líneas o unión de caracteres, aplicados adecuadamente, aseguran que los sistemas atacantes no puedan realizar la separación de caracteres individuales, creando CAPTCHAs seguros.
- Cuidar la implementación, sin sobre-estimar los niveles de seguridad proporcionados por la unión de letras o el uso de líneas, siguiendo las recomendaciones presentadas.
- Crear esquemas alternativos, en función de los nuevos avances en inteligencia artificial, utilizar esos conceptos para renovar el concepto de CAPTCHA, como se comprobaba en el apartado siguiente.

5.2. Características recomendadas para CAPTCHA al mediano plazo

CAPTCHA, como textos relativamente ilegibles está basado en escenarios estáticos donde la respuesta esta visible para todo aquel que pueda resolverla, sin tener que en el proceso mediar ninguna interacción más allá del reconocimiento de patrones lingüístico-numéricos que han sido deformados. No existe ningún tipo de razonamiento ante el estímulo, lo que hace que la nueva generación de CAPTCHAs haga uso de sistemas interactivos para estar un paso delante de la Inteligencia Artificial.

Los siguientes pasos que CAPTCHA debería seguir son:

- Personalizar la experiencia lingüística: Estudios determinaron que a los no angloparlantes tienen dificultades para resolver CAPTCHAs en inglés, por lo que se puede asumir que para los angloparlantes les será de la misma manera difícil resolver CAPTCHAs en otros idiomas que no sean presentados en inglés, por lo que, según la región donde sea desplegado el CAPTCHA, este se presente en un idioma diferente al hablado en el determinado territorio, complicando su resolución en masa.
- Ampliar el tamaño de los textos a ser descifrados: Habitualmente para su resolución se presentan imágenes que contienen dos palabras modificadas. Si se comienzan a añadir oraciones cortas con sentido semántico en lugar de solo palabras al azar, para las personas vendría a ser más fácil la interacción con el sistema al verse ayudados por el idioma y la sintaxis lingüística propia de las personas, aumentando para las maquinas la dificultad al aumentar palabras que segmentar y reconocer.
- Comenzar la creación de CAPTCHAs interactivos: La ventaja que existe entre la brecha humano-maquina es la imposibilidad de las máquinas de razonar frente a estímulos sensitivos interpretados, por lo que CAPTCHA debería evolucionar en un sistema interactivo que requiera obligatoriamente que la respuesta que el usuario de sean fruto de un ejercicio de razonamiento, y que pueden ser de dos tipos.
- De respuesta visible: Se compone de una pregunta de elección en la que el usuario debe escoger la respuesta en función de la solicitud realizada y en donde se muestran las posibles respuestas. En función de la respuesta obtenida, el sistema determinaría la lógica de la misma en función de la pregunta y validaría o no la respuesta. La complejidad del sistema reside en utilizar preguntas que estén en sintonía con la pagina que se esté visitando y que, a pesar de la opción múltiple, la respuesta sea lo suficientemente aparente para los humanos pero que sea un ejercicio de probabilidad para las maquinas.
- De respuesta oculta: Esta prueba juega más con la psicología de las personas, ya que las respuestas dejan de ser aparentes para convertirse en completamente subjetivas en función de la complejidad de pregunta. Las respuestas podrían ser tan sencillas como definir un color en función de parámetros establecidos, hasta manifestación de ideas y sentimientos en frases concretas. Para las personas es un ejercicio fácil de reconocimiento, pero para las maquinas implica encontrar un sentido semántico a la pregunta presentada y darle un contexto en específico antes de emitir una posible respuesta; es decir, la maquina necesita

obligatoriamente razonar y conceptualizar la información que está recibiendo antes de emitir una respuesta válida.

Dado que la neurología aun no es capaz de encontrar una respuesta final a los ejercicios que hacen posible la conciencia y el entendimiento del ser, es que se puede hallar un punto medio entre enseñar a computadoras a leer e interpretar datos fijos y otorgarles la capacidad de razonar sobre lo interpretado y entregar una respuesta fehaciente en correspondencia. A partir de esta diferencia se hace posible que las pruebas sugeridas mantengan la brecha humano-maquina, permitiendo que CAPTCHA pueda seguir siendo usado bajo el mismo concepto como es manejado hoy en día, y mantener la ventaja sobre la Inteligencia Artificial.

6. CONCLUSIONES

Cuando se diseñó Internet para ser usado por el público, no se contempló todo lo que era posible hacer con la red, y todas las formas de innovación que venían con el hecho de poder comunicarnos y compartir información en casi tiempo real. CAPTCHA es uno de esos proyectos que se diseñaron para un objetivo casi específico, pero que el avance de la tecnología hizo que se utilizara en muchos y más variados escenarios, al nivel de convertirse en una prueba de la evolución de la Inteligencia Artificial.

En la actualidad, cuando la vida moderna nos ha hecho casi compartir el espacio físico y virtual con miles de dispositivos electrónicos que están permanentemente conectados a Internet, es cuando se ha vuelto fehaciente la necesidad de siempre diferenciar que es lo que proviene de un ser humano y que proviene de una máquina, considerando como dato que en 2013 más del 58% del tráfico que circulaba por Internet lo producían las maquinas, fruto del Cloud Computing y del nuevo paradigma de información distribuida que se está desarrollando al globalizar la información, el Big Data.

No es realmente posible afirmar cuánto tiempo más podremos disfrutar del relativo control sobre nuestros aparatos electrónicos de la forma en la que veníamos acostumbrados, donde somos nosotros los que utilizamos la tecnología como herramientas de investigación, desarrollo y diversión; para comenzar a considerar a los implementos tecnológicos como seres relativamente sintientes que interactúan con nosotros directamente ayudando y ejecutando tareas sin necesidad de orden alguna. Actualmente estamos en los umbrales de la real inteligencia artificial, la que se aleja de la programación tradicional y se enfoca, fruto de los avances de la neurología y la neurobiología en alimentar a los nuevos sistemas con información sensorial real. Todos estos sistemas en los últimos años dejaron de ser simplemente fruto de programación para convertirse en redes neuronales matriciales que emulan con mucha cercanía el funcionamiento de nuestro cerebro. El presente estudio es una prueba real de cómo estos avances están pisándonos los talones.

La evolución de la Inteligencia Artificial, sumado al descubrimiento de nuevos materiales de construcción de circuitería cada vez más reducida y potente hacen que los avances algorítmicos y físicos comiencen a romper los límites que hace 20 años hubiéramos pensado imposibles para entornos electrónicos, o situaciones solo posibles para el entretenimiento y la ciencia-ficción. Es importante señalar y mantener sobre el tapete que si bien CAPTCHA aún tiene validez como método de seguridad y autenticación, aun con toda la evolución posible para asegurarlo, es un sistema con un tiempo de vida finito y que se acorta a medida que la Inteligencia Artificial avanza a desarrollar algoritmos de conciencia y autonomía.

Parte del desafío de la evolución de CAPTCHA es comenzar a buscar un reemplazo para el mismo que no dependa de la brecha humano-maquina, sino que esté basada en otro tipo de interacciones que efectivamente puedan diferenciar, en un futuro tal vez no tan lejano, el razonamiento humano del razonamiento que pueda ser producido por una entidad autoconsciente. Todos los sistemas que dependan de Internet corren el riesgo de verse superados por el crecimiento y penetración de la tecnología en la sociedad moderna, y es parte del desafío hacia el futuro encontrar soluciones para estos problemas antes que los problemas nos apuren a encontrar las soluciones.

REFERENCIAS

- BBC News, 2013. *Captcha test 'cracked' by US firm vicarious*. Disponible en <http://www.bbc.co.uk/news/technology-24710209>.
- Browning, A., D. Kolas, 2007. *Defeating CAPTCHAs: Applying neural networks*. Virginia Tech. University.
- Bursztein, E., S. Bethard, C. Fabry, J.C. Mitchell, D. Jurafsky, 2010. *How good are humans at solving CAPTCHAs? A Large Scale Evaluation*. Stanford University. Disponible en http://web.stanford.edu/~jurafsky/burszstein_2010_captcha.pdf, 15 pp.
- Bursztein, E., M. Martin, J.C. Mitchell, 2011. *Text-based CAPTCHA strengths and weaknesses*. *ACM Computer and Communication Security*. Disponible en <https://cdn.elie.net/publications/text-based-captcha-strengths-and-weaknesses.pdf>, 14 pp.
- Cai, T., 2008. *CAPTCHA solving with neural networks*. TJHSST Computer Systems Lab 2007-2008.
- Chellapilla, K., P.Y. Simard, 2004. *Using machine learning to break visual human interaction proofs (HIPs)*. NIPS. Disponible en <http://papers.nips.cc/paper/2571-using-machine-learning-to-break-visual-human-interaction-proofs-hips.pdf>, 8 pp.
- Chellapilla, K., K. Larson, P. Simard, M. Czerwinski, 2005. *Computers beat humans at single character recognition in reading based human interaction proofs (HIPs)*. Disponible en <http://ceas.cc/2005/papers/160.pdf>, 8 pp.
- Choi, C.Q., 2013. *Artificial intelligence breaks CAPTCHA protection, coders claim*. Disponible en <http://txchnologist.com/post/65426369724/artificial-intelligence-breaks-captcha-protection>.
- Engber, D., 2014. *Who Made That Captcha?* The New York Times Magazine. Disponible en <http://www.nytimes.com/2014/01/19/magazine/who-made-that-captcha.html>.
- Goodfellow, I.J., Y. Bulatov, J. Ibarz, S. Arnaud, S. Vinay, 2013. *Multi-digit number recognition from Street View Imagery using deep convolutional neural networks*. Cornell University Library. Disponible en <http://arxiv.org/abs/1312.6082>.
- Hawkins, J., G. Dileep, 2011. Hierarchical temporal memory including HTM cortical learning algorithms. White Paper, Numenta.Com.
- Hof, R., 2013. *AI startup vicarious claims milestone* In: Quest To Build A Brain: Cracking CAPTCHA. <http://www.forbes.com/sites/roberthof/2013/10/28/ai-startup-vicarious-claims-milestone-in-quest-to-build-a-brain-cracking-captcha/>.
- Kuchinkas, S., 2013. *Software firm claims breakthrough in computer vision will lead to better AI*. Scientific AmericanTM. Disponible en <http://www.scientificamerican.com/article/ai-captcha-computer-vision/>.
- Naor, M., 1997. *Verification of a human in the loop or Identification via the turing test*. Unpublished Manuscript, Weismann Institute.
- Pérez-Carrasco, J.A., C. Serrano, B. Acha, T. Serrano-Gotarredona, B. Linares-Barranco, *Red neuronal convolucional rápida sin fotogramas para reconocimiento de dígitos*. Disponible en <http://digital.csic.es/bitstream/10261/84753/1/RED%20NEURONAL.pdf>, 4 pp.
- Pinkas, B., T. Sander, 2002. *Securing passwords against dictionary attacks*. In: Proc. ACM Computer and Security Conference (CCS' 02), pp. 161-170. Publisher: ACM Press.
- Pollen, D.A., 2003. *Explicit neural representations, recursive neural networks and conscious visual perception*. Disponible en <http://cercor.oxfordjournals.org/content/13/8/807.long>.
- The CAPTCHA Web Page, 2000. <http://www.captcha.net>. Consultado el 8 de abril de 2014.
- Vicarious FPC Inc., 2012. *Announces \$15 million funding for AI software based on the brain*. Disponible en <http://www.kurzweilai.net/vicarious-announces-15-million-funding-for-ai-software-based-on-the-brain>.
- Vicarious FPC Inc., 2013. *AI passes first Turing Test: CAPTCHA*. Disponible en <http://news.vicarious.com/post/65316134613/vicarious-ai-passes-first-turing-test-captcha>.

- Vicarious FPC Inc., 2013. AI breaks CAPTCHA 'Turing test'. Disponible en <http://www.kurzweilai.net/vicarious-ai-breaks-captcha-turing-test>.
- von Ahn, L., M. Blum, N.J. Hopper, J. Langford, 2003. *CAPTCHA: Using hard AI problems for security*. In: Biham, E. (Ed.). EUROCRYPT 2003, LNCS 2656, pp. 294-311.
- von Ahn, L., M. Blum, J. Langford, 2003. *Telling humans and computers apart (automatically) or how lazy cryptographers do AI*. Disponible en http://www.captcha.net/captcha_cacm.pdf, 4 pp.
- Whitman, R., 2014. *Google's Street View neural network can now decrypt captchas better than a human*. Disponible en <http://www.extremetech.com/computing/174275-google-has-built-a-neural-network-to-identify-100-million-house-numbers-for-streetview>.