

UMA METODOLOGIA DE ATRIBUIÇÃO DE AUTORIA APLICADA A INVESTIGAÇÕES SOBRE ABUSO SEXUAL INFANTIL

AN AUTHORSHIP ATTRIBUTION MODEL APPLIED TO PEDOPHILIA CRIME INVESTIGATIONS

Aurélio Julbert de Assis Ruprecht^a

Marcelo da Silva Moreira^b

Enrique Muriel Torrado^c

Moisés Lima Dutra^d

RESUMO

Objetivos: Identificar o atual estado da arte das pesquisas científicas no campo da atribuição de autoria aplicada a investigações de crimes sexuais contra crianças e adolescentes pela Internet envolvendo material escrito. Propor uma metodologia de utilização da atribuição de autoria para identificação de suspeitos de serem autores de textos com conteúdo que incentive o abuso sexual infantojuvenil. **Metodologia:** Trata-se de uma pesquisa qualitativa que utiliza a Revisão Sistemática da Literatura para identificar trabalhos que versem a respeito das técnicas de atribuição de autoria com o intuito de se buscar evidências científicas de sua aplicação a problemas semelhantes ao abordado no presente estudo. **Resultados:** Apresenta-se o atual estado da arte das pesquisas científicas que relacionam a utilização de técnicas de atribuição de autoria a textos presentes na internet que incentivam a prática de abuso sexual de crianças e adolescentes e, a partir disso, propõe-se uma metodologia para identificação de autores de textos com aquelas características. **Conclusões:** Conclui-se que não existe abundância de pesquisas científicas sobre esse tema, o que sugere ser um campo aberto à novos estudos. Também se conclui que é plenamente possível a aplicação das técnicas de atribuição de autoria na identificação dos prováveis autores de textos que tenham como objetivo orientar e fomentar a prática de abuso sexual infantojuvenil, o que é explicitado pela metodologia proposta.

^a Mestre em Ciência da Informação pela Universidade Federal de Santa Catarina (UFSC). Escrivão da Polícia Federal lotado na Diretoria de Combate ao Crime Organizado. Direção-Geral da Polícia Federal. E-mail: ajar1971@gmail.com

^b Mestre em Ciência da Informação pela Universidade Federal de Santa Catarina (UFSC). Perito Criminal da Polícia Federal. E-mail: marcelomoreira3000@gmail.com

^c Doutor em Información Científica pela Universidad de Granada. Docente no Programa de Pós-Graduação em Ciência da Informação da Universidade Federal de Santa Catarina (UFSC). E-mail: enrique.muriel@ufsc.br

^d Doutor em Computação pela Universidade de Lyon, França. Docente do departamento de Ciência da Informação da Universidade Federal de Santa Catarina (UFSC). E-mail: moises.cin.ufsc@gmail.com

Descritores: Abuso sexual infantojuvenil pela internet. Atribuição de autoria. Estilometria. Pedofilia. Investigação Policial.

1 INTRODUÇÃO

As indiscutíveis facilidades oferecidas pela Internet têm sido usadas também para a prática de diferentes crimes, como estelionato, furto de ativos bancários, acessos indevidos a informações de dispositivos pessoais, comerciais e públicos, mas também para a prática de crimes relacionados ao abuso sexual infantil e de adolescentes, comumente referido pelo termo pedofilia. O cometimento de delitos usando a rede mundial de computadores é incentivado pela dificuldade de se identificar seus autores (ABBASI; CHEN, 2006; ZHENG *et al.*, 2006; BHARGAVA; MEHNDIRATTA; ASAWA, 2013; YANG; CHOW, 2014). Os crimes relacionados ao abuso sexual de menores de idade pela Internet são praticados por meio de diversas modalidades, incluindo a atração e cooptação de crianças e adolescentes e a troca, comércio e divulgação de fotos e filmes digitais de práticas pedófilas ou de pornografia infantil (FRANCO; MAGALHÃES, 2015).

Um dos problemas frequentemente associados ao anonimato online é que isso dificulta a responsabilidade social, como comprovado pelos altos níveis de cibercrime. Embora as pistas de identidade sejam escassas no ciberespaço, os indivíduos geralmente deixam rastros de identidade textuais (ABBASI; CHEN, 2008, p. 1, tradução nossa).

O Estatuto da Criança e do Adolescente (Lei n.º 8.069, de 13 de julho de 1990) prescreve um conjunto de normas que tipificam como crime diversas práticas que podem causar algum dano físico ou psicológico a este grupo vulnerável. Recentes investigações conduzidas pela Polícia Federal, no Brasil, identificaram o uso de documentos com características bem peculiares que funcionam como incentivo a essa prática criminosa. Trata-se de verdadeiras apostilas que ensinam como praticar sexo com menores de idade (MOREIRA, 2020). Esses documentos, que podem ser chamados de “manuais de pedofilia”, são guias didaticamente organizados para “ensinar” a um adulto a prática de “sexo seguro” com crianças e adolescentes. Por “sexo seguro”, os “manuais” se referem à segurança do ponto de vista da saúde dos adultos e das vítimas, mas

também se referem a evitar que aqueles que o praticam sejam identificados e presos. Tais “manuais de pedofilia” incluem informações detalhadas sobre sexualidade infantil, tabus, riscos legais, técnicas de como encontrar crianças, como seduzir as vítimas e de como se aproximar dos seus responsáveis e conquistar sua confiança, entre outras.

De fato, esse tipo de material é elaborado de forma a ser autoexplicativo e de fácil leitura. São organizados em capítulos encadeados de forma didática, contendo ilustrações e repletos de exemplos práticos, tem, portanto, o objetivo de incentivar novos adeptos à prática de sexo com menores de idade, perpetuando essa atividade criminosa (MOREIRA, 2020).

Evidentemente, esses tipos de texto não trazem a discriminação de sua autoria, apenas codinomes, e são divulgados da forma mais discreta possível, normalmente usando camadas da Internet de difícil acesso, como a *Dark Web*, em comunidades privadas de acesso restrito aos seus membros, que se dedicam a ensinar teorias, métodos e técnicas para consumação de abuso sexual contra crianças e adolescentes.

Assim, a importância de identificar os autores desses guias é evidenciada pela necessidade de colocar à disposição da Justiça aqueles que cometem o crime de autoria dos manuais, já que a sua simples publicação é considerada delito de incitação ao crime, conforme o Art. 286 do Código Penal Brasileiro (BRASIL, 1940).

O que o presente artigo propõe é a identificação de indícios da autoria de manuais de pedofilia por meio da análise do estilo de redação empregado, através de ferramentas de análise do estilo de escrita (estilometria) e de atribuição de autoria. Parte-se da premissa de que cada indivíduo tem um estilo próprio de escrita, de forma que, dentro de uma margem de erro aceitável, é possível identificar o autor de um texto pela forma como ele foi escrito, ou pelo menos elencar um rol restrito de possíveis suspeitos. Características como vocabulário escolhido e a maneira de usá-lo, emprego de verbos, semântica, concordância verbal e nominal, organização do texto, entre outras, combinadas com dados técnicos do arquivo eletrônico (metadados), funcionariam como a

assinatura de seu autor, da mesma forma que se pode reconhecer um compositor ou pintor pela análise de uma peça de sua produção artística.

Classicamente, a atribuição de autoria parte da comparação entre as características de escrita presentes em obras de autores conhecidos e as de textos cuja autoria pretende-se conhecer. Essa comparação permite associar estilos de escrita semelhantes de forma a atribuir ao texto de autoria desconhecida um provável autor. Mas, e se o problema for identificar a autoria de um texto sem que se tenha um banco de dados contendo perfis estilométricos de textos com autores conhecidos a serem comparados? Essa é a questão central desta pesquisa: como identificar a autoria de textos com conteúdo de incentivo à prática de crimes de abuso sexual de crianças e adolescentes, em formato eletrônico e encontrados na Internet, sem que se tenha amostras de textos de autoria conhecida para a comparação? Esse é exatamente o problema encontrado na vida prática, o de identificar os autores dos “manuais de pedofilia” sem que se tenha materiais textuais produzidos por suspeitos para a comparação, ou mesmo sem que sequer haja suspeitos.

O objetivo geral deste estudo é propor uma metodologia para a identificação dos autores dos “manuais de pedofilia” com base em seu estilo de escrita, através da utilização das técnicas de estilometria e atribuição de autoria. Especificamente, buscou-se trazer para o presente estudo os pontos basilares das técnicas de atribuição de autoria e estilometria colhidos da literatura científica sobre o tema. Secundariamente, objetivou-se também identificar o atual estado da arte das pesquisas científicas no campo da atribuição de autoria aplicada a investigações sobre autores de crimes sexuais contra crianças e adolescentes pela Internet.

Este trabalho é composto pelas seguintes seções, além da introdução: uma segunda seção que mostra o panorama dos casos denunciados no Brasil de abuso sexual de crianças e adolescentes pela Internet ou não. Pela terceira seção sobre a metodologia de pesquisa empregada nesse trabalho. A quarta seção traz a análise dos trabalhos científicos que contribuíram para a discussão e se divide em uma subseção com os artigos relativos aos conceitos basilares da estilometria e atribuição de autoria e outra subseção contendo trabalhos

correlatos. A quinta seção apresenta os resultados da pesquisa, incluindo subseção contendo a proposta de metodologia para a aplicação das técnicas de atribuição de autoria na identificação de material objeto desse estudo. A sexta seção traz as conclusões e perspectivas futuras.

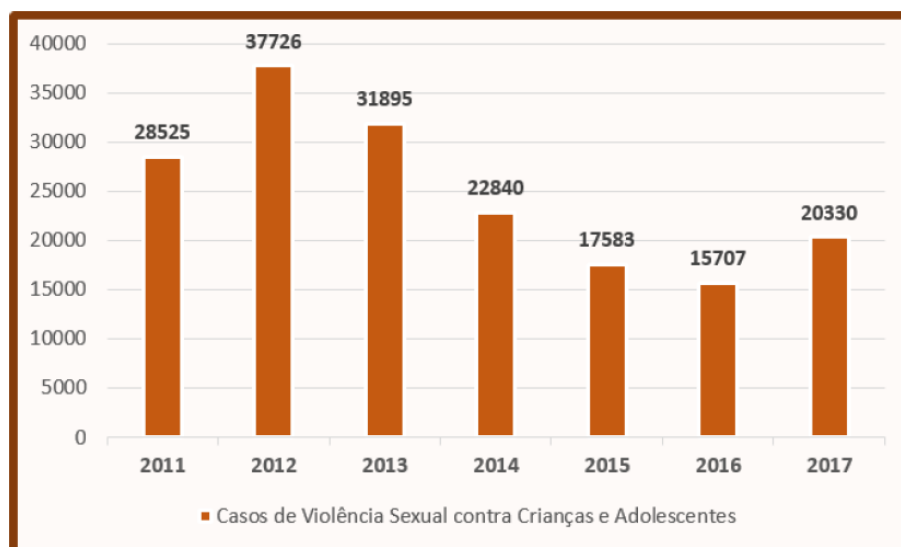
2 ABUSO SEXUAL DE CRIANÇAS E ADOLESCENTES NO BRASIL

Esta seção mostra um panorama do problema do abuso sexual de crianças e adolescentes no Brasil em anos recentes a partir de dados de denúncias feitas às Organizações da Sociedade Civil (OSC) e ao Ministério da Mulher, da Família e dos Direitos Humanos (MMFDH), que em âmbito nacional é o órgão vinculado à Presidência da República responsável pela coordenação das ações de combate dessas violações.

Iniciando-se pelos dados coletados junto ao MMFDH, tem-se os números das comunicações de violações contra crianças e adolescentes feitas pelo Disque 100, canal oficial de denúncias do referido Ministério, que revelaram a evolução histórica dos casos de violência sexual infantojuvenil no Brasil de 2011 a 2017.

Os valores absolutos de casos de violência sexual são expressivos tendo, em 2012 atingido seu maior valor: 37.726 denúncias (MINISTÉRIO DA MULHER, DA FAMÍLIA E DOS DIREITOS HUMANOS, 2018 *apud* MOREIRA, 2020).

Gráfico 1 - Histórico de casos de violência sexual contra crianças e adolescentes no Disque 100

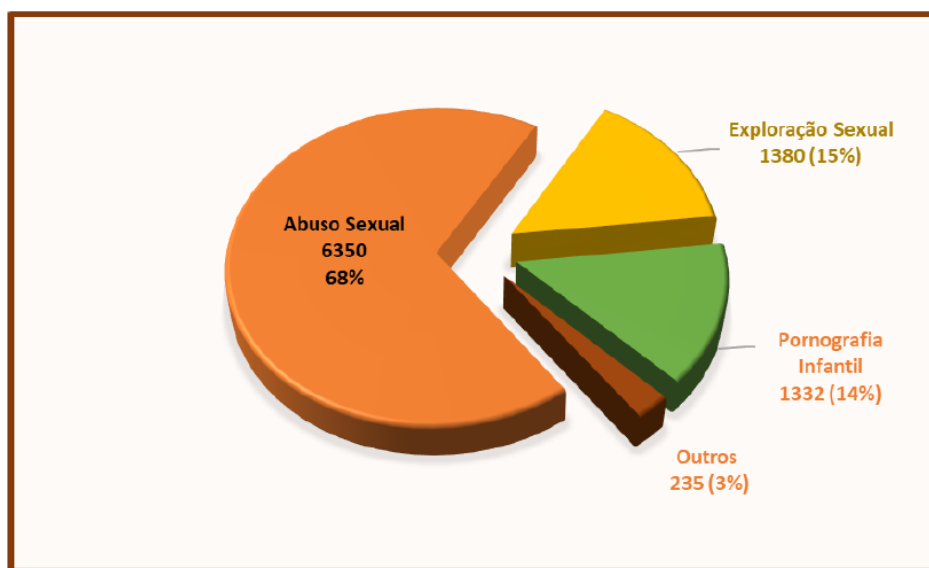


Fonte: Adaptado de MMFDH (2018 *apud* MOREIRA, 2020).

Ao observar com mais detalhes os dados mais recentes do MMFDH, referentes ao primeiro semestre de 2018, percebe-se que o Disque 100 recebeu 36.757 denúncias distribuídas entre todos os tipos de violações contra crianças e adolescentes no Brasil, a negligência, com 26.242 casos, é a mais frequente violação, seguida da violência psicológica (17.031 casos), da violência física (14.355) e da violência sexual (9.297 casos).

Assim, nota-se que violência sexual é a quarta violação mais denunciada no Disque 100, sendo que 68% desses casos estão relacionados a situações de abuso sexual infantil; 15% dos casos, à exploração sexual de menores de idade, 14% dos casos, à pornografia infantil, e, 3% a outras situações. O Gráfico 2 retrata a forma como estão distribuídas as 9.297 denúncias de violência sexual contra crianças e adolescentes feitas no primeiro semestre de 2018.

Gráfico 2 - Principais denúncias de violência sexual contra menores de idade no Disque 100 (janeiro a junho de 2018)



Fonte: Adaptado de MMFDH (2018 *apud* MOREIRA, 2020).

Os dados mais atuais evidenciam que o problema se agravou. No primeiro semestre de 2020 o número de denúncias ao Disque 100 cresceu intensamente, chegando a mais de 53 mil denúncias de violências cometidas contra crianças e adolescentes, dentre as quais 9 mil foram de cunho sexual, e no segundo semestre do mesmo ano o número de denúncias atingiu 36 mil, perfazendo uma marca anual de 89 mil denúncias de violência contra crianças e adolescentes.

É importante lembrar que no primeiro semestre de 2020 o País viveu a primeira onda da pandemia de Covid-19, provocando o isolamento social, suspensão de aulas, desemprego e fechamento de empresas, induzindo o acirramento da crise econômica, elementos que podem ter contribuído para o aumento. Apenas entre janeiro e maio de 2021, o Disque 100 recebeu mais de 35 mil denúncias de violência contra crianças e adolescentes, das quais 17,5%, ou mais de 6 mil, foram denúncias de violência sexual.

A Childhood Brasil é uma organização da sociedade civil (OSC) que tem como objetivo a proteção à infância e à adolescência com foco no enfrentamento do abuso e da exploração sexual contra esse público (CHILDHOOD BRASIL, 2021a). A Childhood Brasil apurou que 70% das pessoas estupradas no Brasil anualmente são crianças e adolescentes (CHILDHOOD BRASIL, 2021b).

O quadro da violência sexual contra crianças e adolescentes, quando cometida por meio da Internet, pode ser medido pelas estatísticas levantadas pela OSC Safernet. A Safernet é uma OSC dedicada à promoção e defesa dos Direitos Humanos na Internet no Brasil” (SAFERNET, 2021). Essa OSC criou e mantém uma central nacional de denúncias de crimes cibernéticos em parceria com Ministérios Públicos e com a Secretaria de Direitos Humanos da Presidência da República (SDH), tornando-se referência nacional no combate a crimes cibernéticos, notadamente os de abuso sexual contra crianças e adolescentes (pornografia infantil).

De acordo com a Safernet, as denúncias de pornografia infantil bateram o recorde em 2020, com 8 mil denúncias, mais do que o dobro do ano anterior. Contudo, os números mais recentes indicam o crescimento de 33,45% nos quatro primeiros meses de 2021 em comparação com o mesmo período de 2020.

3 METODOLOGIA DA PESQUISA

A metodologia de pesquisa adotada foi a revisão sistemática da literatura (RSL) com o intuito de se buscar evidências científicas da aplicação das técnicas de atribuição de autoria a problemas semelhantes ao do presente estudo.

A RSL foi executada no dia 27/11/2019 e usou como termos de busca (*paedophilia OR pedophilia OR paedophile OR pedophile*) AND "*authorship attribution*" e "*authorship attribution*", portanto buscou apenas textos no idioma inglês. As bases de dados utilizadas foram: a IEEE, ACM, COMPENDEX, Science Direct, SCOPUS e Web of Science. Os resultados das buscas são relatados mais à frente nesse texto.

Os textos selecionados trouxeram o conhecimento teórico a respeito das técnicas de estilometria e atribuição de autoria que fundamentaram a metodologia proposta por esse estudo.

4 ESTILOMETRIA E ATRIBUIÇÃO DE AUTORIA

A estilometria e a atribuição de autoria estão intimamente relacionadas. A primeira tem origem na linguística e pode ser entendida como o “estudo estatístico do conjunto de características textuais (características estilométricas) da obra de um autor” (ZHENG *et al.*, 2006, p. 379, tradução nossa). Com relação à atribuição de autoria, esta trata de tentar identificar o autor de um texto, anônimo ou não, de cuja autoria se tem dúvida (PENG *et al.*, 2003). A atribuição de autoria busca, através de diferentes técnicas — inclusive computacionais —, descobrir a autoria de um texto pela comparação entre suas características estilométricas e o conjunto das características de escrita de autores conhecidos (ABBASI; CHEN, 2006). Em outras palavras, a atribuição de autoria trata de identificar autores a partir de seus textos com base nas suas características únicas de escrita (GE; SUN; SMITH, 2016; RAMNIAL; PANCHOO; PUDARUTH, 2015).

A questão tradicional que a atribuição de autoria se propõe a resolver pode ser explicada pela definição de Zheng *et al.* (2006), a qual afirma que, dado um conjunto de textos de vários autores, deve-se atribuir a autoria de um novo texto para um deles (ZHENG *et al.*, 2006; STAMATATOS, 2009). Trata-se, portanto, de um teste de hipótese estatística ou um problema de classificação, sendo que a “essência dessa classificação é identificar um conjunto de recursos

que permanecem relativamente constantes em um grande número de escritos criados pela mesma pessoa” (ZHENG *et al.*, 2006, p. 380, tradução nossa).

Nos últimos anos a atribuição de autoria de mensagens anônimas recebeu atenção notável nas comunidades de análise forense e de mineração de dados. Durante o as últimas duas décadas essa técnica se estendeu à comunicação facilitada por computador ou a documentos on-line (como e-mails, SMS, tweets, mensagens de chat instantâneas etc.) para levar terroristas, pedófilos e golpistas à Justiça (BHARGAVA; MEHNDIRATTA; ASAWA, 2013, p. 37, tradução nossa) .

Um exemplo conhecido de atribuição de autoria com base na estilometria é o caso dos escritos federalistas (*Federalist Papers*), um conjunto de 146 ensaios sobre política escritos por John Jay, Alexander Hamilton e James Madison. Doze desses ensaios tinham a autoria disputada entre Hamilton e Madison e a aplicação da análise do estilo de escrita dos autores envolvidos foi fundamental para se atribuir a autoria correta ao subconjunto de textos em dúvida (PENG *et al.*, 2003; ZHENG *et al.*, 2006; STAMATATOS, 2009).

As características estilométricas se referem à maneira particular de um autor usar recursos gramaticais em seus textos de modo a marcar seu estilo de escrita. Podem ser características lexicais, sintáticas, estruturais, de conteúdo, relativas à extensão das palavras usadas e ainda características idiossincráticas (BHARGAVA; MEHNDIRATTA; ASAWA, 2013, p. 37, tradução nossa). O aspecto lexical se refere à possibilidade de medição estatística da variação de ocorrências de palavras, denunciando o vocabulário de um autor. As características sintáticas incluem a função das palavras no texto, pontuação e marcadores de discurso. Palavras funcionais são aquelas que expressam uma relação gramatical ou estrutural com outras palavras em uma mesma sentença, e são úteis para discriminar a autoria porque a variação do uso dessas palavras reflete fortemente as escolhas estilísticas.

Os aspectos estruturais do texto são especialmente úteis para análise de textos online, podendo incluir atributos de organização do texto, extensão das palavras ou sentenças e o layout típico dos textos de um autor. Outras características técnicas podem ser analisadas ao se tratar de textos digitais, como a variação do tamanho do arquivo eletrônico, metadados e as fontes usadas, seu tamanho e cores. O conteúdo específico do texto pode revelar o uso

de palavras-chave específicas, chavões, terminologia técnica ou frases comuns a certos tópicos ou atividades profissionais específicas. Os atributos idiossincráticos da escrita de um autor se referem a um modo peculiar de se manifestar, que pode ser reconhecido em erros de ortografia, erros gramaticais e outras anomalias recorrentes (ABBASSI; CHEN, 2008). Como as características de escrita podem variar enormemente, “mais de 1000 atributos diferentes têm sido usados em análises de autoria de textos sem que um consenso tenha sido alcançado sobre a melhor combinação desses atributos” (RUDMAN, 1997 *apud* ABBASSI; CHEN, 2008, p. 4, tradução nossa). Porém, aceita-se de maneira geral que o desempenho da análise de autoria depende da combinação entre as características de texto selecionadas e técnicas de análise (ZHENG *et al.*, 2006).

As técnicas de análise de estilometria podem ser categorizadas em supervisionadas e não supervisionadas, com abordagem estatística ou baseada em aprendizagem de máquina. Técnicas supervisionadas são aquelas que exigem a identificação das classes de autor para categorização, enquanto as técnicas não supervisionadas realizam a categorização sem se conhecer as classes de autores. As técnicas supervisionadas incluem *support vector machines* (SVM), redes neurais, árvores de decisão e análise discriminante linear. As abordagens estatísticas e de aprendizagem de máquina são as mais usadas, sendo que as técnicas estatísticas têm o benefício de fornecer maior potencial explicativo. Isso pode ser útil para avaliar tendências e variações em quantidades maiores de texto, mas o aumento da capacidade computacional disponível fez surgir técnicas de aprendizado de máquina como as acima citadas (ABBASSI; CHEN, 2006).

Os métodos de atribuição de autoria podem ainda ser classificados em baseados em perfis de autores ou baseados em amostras de textos. Para a abordagem baseada no perfil do autor, uma única representação ou perfil estilométrico é produzido para cada autor usando dados de treinamento, e então cada texto de autoria desconhecida é comparado com os perfis de modo a se poder atribuir ao texto um provável autor. A abordagem baseada em perfil é capaz de lidar com textos muito curtos, pois concatena todos os textos do mesmo

autor (YANG; CHOW, 2014). Quando o método de atribuição de autoria segue a abordagem baseada em amostras, é produzida uma representação do estilo de escrita do autor por texto de treinamento, além de um modelo de classificação que é construído para estimar o autor mais provável de um texto anônimo. Comparado com o modelo de abordagem baseada em perfil, a abordagem por amostra combina mais facilmente diferentes características de texto e é mais robusta quando o tamanho do conjunto de autores candidatos é amplo (YANG; CHOW, 2014; STAMATATOS, 2009; KOURTIS; STAMATATOS, 2011). Um estudo conduzido por Kourtis e Stamatatos (2011) combinou duas técnicas de atribuição de autoria que representam bem cada uma das abordagens acima. A conclusão dos autores foi que as abordagens baseadas em perfil e as fundadas em amostragens podem ser complementares entre si e que aplicar as duas abordagens em conjunto pode melhorar significativamente o desempenho de uma atribuição de autoria (KOURTIS; STAMATATOS, 2011). A atribuição de autoria pode enfrentar obstáculos como a questão do tamanho das amostras de texto disponíveis, seja para o treinamento dos algoritmos ou para fins de comparação entre os escritos de autoria conhecida e os de autoria desconhecida a ser determinada. Essa questão impacta diretamente no grau de precisão dos métodos de atribuição de autoria (ZHENG *et al.*, 2006; STAMATATOS, 2009; RAMNIAL; PANCHOO; PUDARUTH, 2015; ROCHA *et al.*, 2017).

Como visto até aqui, a atribuição de autoria está fundamentada na estilometria. O funcionamento do processo de atribuição de autoria clássico é composto por duas tarefas: i) a identificação das características únicas de escrita de diversos autores com o objetivo de criar um banco de dados de estilos para futuras comparações; e ii) a análise do estilo de escrita usado num texto de autoria desconhecida ou duvidosa. Para se chegar à atribuição de autoria a um texto, é necessário efetuar uma comparação entre os resultados dessas duas tarefas.

4.1 TRABALHOS CORRELATOS

A Tabela 1 sintetiza o resultado da revisão sistemática de literatura (RSL). A primeira coluna relaciona às bases de dados nas quais a RSL foi executada e as outras colunas mostram o número de trabalhos recuperados de acordo com o termo de busca aplicado.

A intenção da comparação foi mostrar o espaço que os estudos dedicados a aplicar as técnicas de atribuição de autoria ao problema de identificação de pedófilos ocupam no universo de trabalhos científicos sobre o tema atribuição de autoria. Como se pode perceber da comparação proposta, aparentemente não existe a preocupação específica de se usar a atribuição de autoria como forma de identificar autores de textos pedófilos ou que incentivem o abuso sexual de crianças e adolescentes pela Internet.

Tabela 1 - Resultados da Revisão Sistemática da Literatura

Termo de busca/Base	(paedophilia OR pedophilia OR paedophile OR pedophile) AND "authorship attribution"	"authorship attribution"
IEEE	0	1.279
ACM	0	65
COMPENDEX	1	543
SCIENCE DIRECT	4	258
SCOPUS	1 (mesmo trabalho do Compendex)	831
WEB OF SCIENCE	0	566

Fonte: Dados da pesquisa (2019).

A análise dos resultados da RSL e dos textos que ela retornou mostram que, além de representarem um número relativamente pequeno, os trabalhos que correlacionam atribuição de autoria à pedofilia também não são especificamente direcionados ao tema. Na maioria das vezes, apenas citam que as técnicas poderiam ser usadas na identificação de autores de crimes de pedofilia, entre outras aplicações. Para ilustrar, toma-se como exemplo o artigo de Ishihara (2014), extraído da base Compendex, que relata o cotejo entre dois diferentes procedimentos forenses de comparação de textos, usando "mensagens de *chatlog*, que foram apresentadas como evidência para processar pedófilos para teste" (ISHIHARA, 2014, tradução nossa), contudo os métodos em estudo já eram preexistentes e as amostras de mensagens foram usadas

unicamente como elementos de teste. O estudo teve o objetivo comparar duas técnicas de cálculo de atribuição de autoria com base nas taxas de verossimilhança.

O Quadro 1 evidencia a contribuição dos artigos selecionados quanto à relação do estudo com os crimes de abuso sexual de crianças e adolescentes cometidos pela Internet.

Quadro 1 - Relação entre o estudo selecionado e crimes de abuso sexual de crianças e adolescentes cometidos pela internet

Título	Autores	Base	Relação entre o estudo e pedofilia pela internet
A Comparative Study of Likelihood Ratio Based Forensic Text Comparison Procedures: Multivariate Kernel Density with Lexical Features vs. Word N-grams vs. Character N-grams	Shunichi Ishihara	Compendex	Usa mensagens de <i>chatlog</i> , que foram apresentadas como evidência para processar pedófilos para teste
A feature selection method for author identification in interactive communications based on supervise learning and language typicality	Esther Villar-Rodriguez; Javier Del Ser; Miren Nekane Bilbao; Sancho Salcedo-Sanz.	Science Direct	Sugere uma lista de aplicações possíveis para o método proposto pelo estudo, citando que pode ser usado para identificação de pedófilos, contudo, o método não foi desenvolvido para este propósito de modo específico.
Early detection of deception and aggressiveness using profile-based representations	Hugo Jair Escalante; Esaú Villatoro-Tello; Sara, E. Garza; A. Pastor López-Monroy; Manuel Montes-y-Gómez; Luis Villa señor-Pineda	Science Direct	Esse estudo é focado na proposição de uma metodologia para detectar ataques potenciais do tipo predadores sexuais que se aproximam de menores ou de usuários agressivos, o mais cedo possível.
Strength of linguistic text evidence: A fused forensic text comparison system	Shunichi Ishihara	Science Direct	Usa mensagens de <i>chatlog</i> , que foram apresentadas como evidência para processar pedófilos para teste
Towards an integrated e-mail forensic analysis framework	Rachid Hadjidj; Mourad Debbabi; Hakim Lounis; Farkhund Iqbal; Adam Szporer; Djamel Benredjem	Science Direct	O trabalho foca os crimes cibernéticos que se utilizam de mensagens por e-mail, cita os crimes de abuso sexual de crianças e adolescentes, nominando-os por pedofilia, cita que o anonimato na comunicação por e-mail é um dos principais artifícios explorado

			por terroristas, pedófilos e golpistas.
--	--	--	---

Fonte: Elaborado pelos autores (2019).

Os resultados da execução da RSL justificam plenamente mais pesquisas sobre a aplicação de técnicas já conhecidas de atribuição de autoria na identificação dos autores de crimes de abuso infantojuvenil, além de justificarem também o desenvolvimento de novas abordagens para fins forenses.

5 RESULTADOS

O objetivo principal deste estudo é colaborar para a solução de um problema social dos tempos atuais, os crimes de abuso sexual de crianças e adolescentes pela Internet, especificamente quando praticados por meio da disseminação de textos nomeados neste trabalho como “manuais de pedofilia”. Diante da natureza e característica desse tipo de material, seu formato de texto, adotou-se desde o início da pesquisa a busca na produção científica por soluções a partir das técnicas de estilometria e atribuição de autoria. Esse esforço produziu dois resultados. O principal é o atingimento do objetivo proposto através da elaboração de uma metodologia de aplicação das técnicas já conhecidas de atribuição de autoria conjugadas com as de estilometria, mas orientada originariamente para a identificação dos autores dos chamados “manuais de pedofilia”. Um outro resultado acessório da pesquisa foi a constatação de uma lacuna na produção científica sobre a relação dos temas “estilometria” e “atribuição de autoria” com o tema “pedofilia”. Para o levantamento da fundamentação teórica deste trabalho usou-se a metodologia da RSL, tendo como conclusão a percepção da baixa produção científica nessas áreas.

5.1 A PROPOSTA DE METODOLOGIA PARA APLICAÇÃO DE ATRIBUIÇÃO DE AUTORIA NA IDENTIFICAÇÃO DE TEXTOS DO TIPO “MANUAIS DE PEDOFILIA”

A solução que se propõe para a questão central deste artigo começa com a inversão da ordem pela qual comumente as técnicas de atribuição de autoria

funcionam. Como visto anteriormente, a técnica de atribuição de autoria mais aplicada costuma partir da análise estilométrica de textos de autores conhecidos, formando para cada autor um conjunto próprio de características de estilo de escrita. Desse ponto em diante aplica-se o mesmo tipo de análise sobre os textos de autoria desconhecida para classificá-los de acordo com suas características identitárias. Na sequência, a forma mais usual de utilização da atribuição de autoria faz a comparação entre as características estilométricas dos textos cuja autoria se quer conhecer e o conjunto das características de estilo de vários autores já registrados, resultando na indicação probabilística quanto à autoria do texto alvo da análise.

De acordo com a metodologia aqui proposta, a primeira fase é a identificação de textos suspeitos pelas forças policiais. É a Fase 1 da Figura 1. Nessa primeira fase fica evidente a importância da participação do policial no processo, pois é por meio de seu trabalho minucioso de combate aos crimes popularmente identificados como “pedofilia” que são descobertos os indícios materiais dos mesmos, como, por exemplo, os “manuais de pedofilia” objetos deste estudo. Essa fase inicial se dá em momento anterior ao processamento computacional e linguístico dos textos.

A etapa seguinte (Fase 2) do modelo ora proposto é a execução da análise estilométrica de textos cujo conteúdo se encaixa na descrição do que foi chamado neste estudo de “manuais de pedofilia” ou na descrição de outros tipos de texto que apresentem relação com crimes de abuso sexual de crianças e adolescentes descobertos por investigações policiais. Essa fase vai informar todo o processo e é executada por algoritmos dedicados, de domínio público ou desenvolvidos para esse fim específico, podendo-se alternar os algoritmos usados de acordo com características do material em análise. Nesta segunda fase os textos suspeitos são examinados até a formação de sua identidade estilométrica. Com a análise estilométrica de textos de autoria desconhecida disponíveis, tem-se o perfil estilométrico do autor de cada texto analisado. Essa informação dá início à Fase 3, subdividida em fases 3a e 3b.

O que ocorre na Fase 3a é a classificação dos textos suspeitos de acordo com suas características estilométricas próprias. Nesse estágio do processo os

textos de autoria desconhecida são agrupados de acordo com o seu perfil estilométrico mapeado na fase anterior, formando subconjuntos de texto (na Figura 1, são os subconjuntos 1, 2, n) que se diferenciam uns dos outros de forma inequívoca. Aí torna-se possível identificar os traços de seus autores, mesmo que eles ainda sejam desconhecidos. Dessa forma, ao longo do tempo esses subconjuntos passam a ser ampliados com a adição de novos textos que tenham perfil estilométrico de escrita semelhante, permitindo processos de análise estilométrica cada vez mais confiáveis e robustos. Como ato contínuo, ainda na Fase 3a, cada subconjunto é “batizado” por um código identificador único que será posteriormente substituído pelo nome do autor, quando este for descoberto. A fase 3b é constituída pela formação de um banco de dados específico que alimentará, na fase seguinte, o algoritmo responsável pela busca por estilos de escrita idênticos em materiais de autoria conhecida.

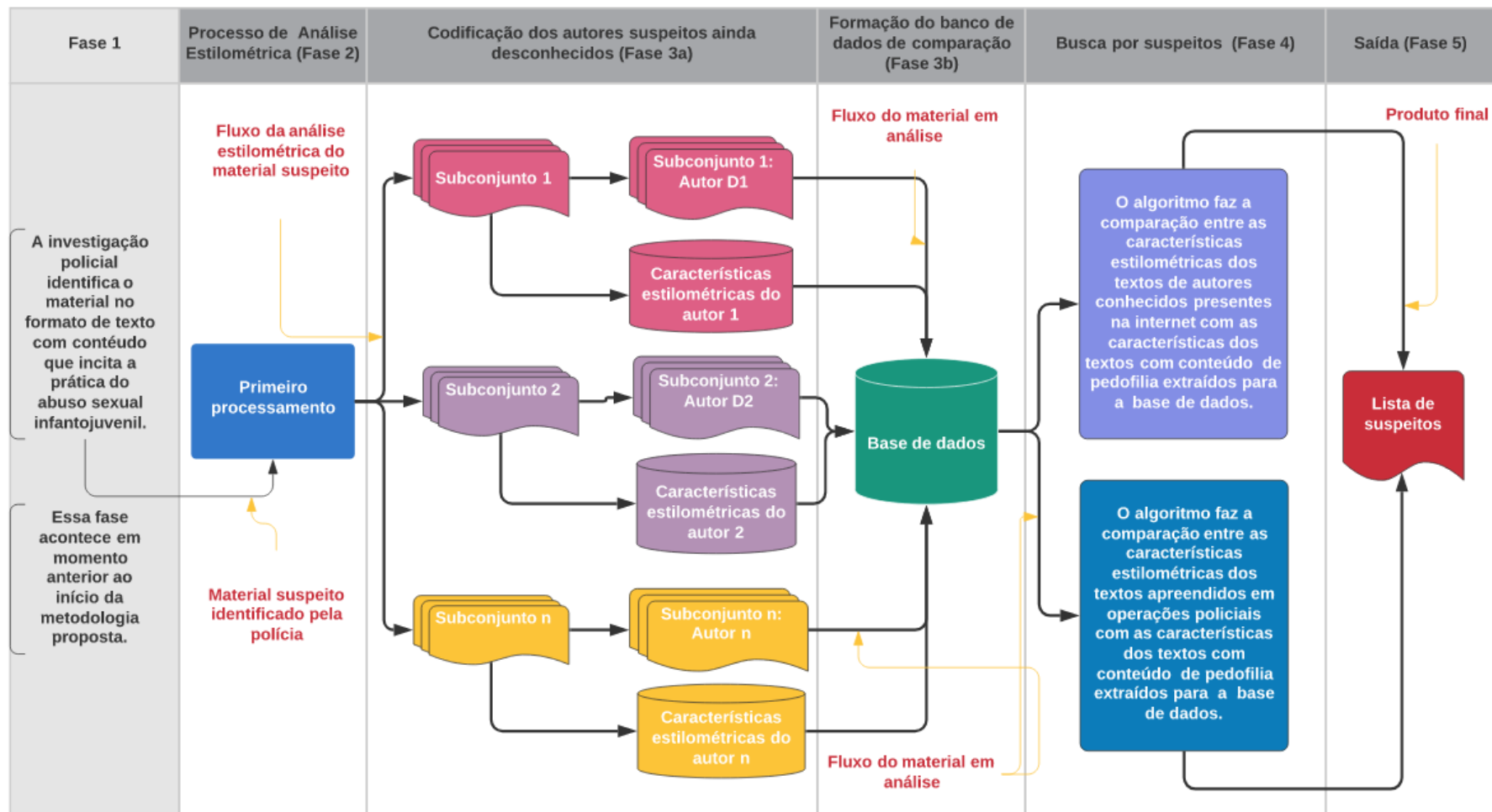
O momento seguinte (Fase 4) do processo pode ser descrito como sendo um verdadeiro trabalho de investigação policial, de forma autônoma e ininterrupta, no qual algoritmos especialmente desenvolvidos para esse fim processam, em diferentes ambientes na Internet, textos de autoria conhecida que comumente são encontrados em: i) sites pessoais; ii) perfis em redes sociais; iii) sites profissionais, tais como de empresas ou aqueles em que o profissional explana seus produtos e serviços; iv) perfis em redes de relacionamento voltadas para a profissão, como o LinkedIn; v) fóruns de discussão profissionais, como aqueles onde profissionais de uma determinada área trocam conhecimentos, além de fazer e responder perguntas; vi) bancos de dados de publicações de diferentes áreas, tais como periódicos; entre outros. De maneira geral, são locais na Internet escolhidos pelo critério da maior possibilidade de se coletar grandes volumes de texto de um mesmo autor. O que esse algoritmo faz é a comparação entre as características dos textos de autores conhecidos presentes nesses ambientes da Internet e as características dos textos cujo conteúdo seja suspeito de conter material criminoso, extraídos para a base de dados formada na Fase 3b. O algoritmo é previamente “treinado” para a realização do trabalho de atribuição de autoria (comparação) e provido dos dados estilométricos dos autores dos textos pedofílicos, já codificados.

Da comparação feita pelo algoritmo se tem como resultado uma lista de suspeitos (Fase 5), que será entregue aos agentes que investigam esses crimes para ser usada como informação de inteligência e para outras providências da investigação.

Outra possibilidade de aplicação do algoritmo seria sobre textos produzidos por suspeitos em investigação, apreendidos em operações policiais ou não. Ou seja, uma investigação policial toma posse de textos produzidos por um suspeito de cometimento de crimes de abuso sexual de crianças ou adolescentes, mas que tem a identidade conhecida, e compara o material apreendido com as informações estilométricas mapeadas nas Fases 2 e 3a e armazenadas em banco de dados específico na fase 3b do modelo proposto. A finalidade de tal procedimento é a descoberta de materiais probatórios que sirvam à confirmação ou negação quanto ao envolvimento anterior do suspeito nesse tipo de crime.

Em ambas as formas de aplicação do método não se chegaria a uma conclusão inequívoca sobre a autoria de textos pedofílicos, mas certamente se reduziria a lista de suspeitos. Esta metodologia, usada em conjunto com outras técnicas de investigação, poderia ajudar na identificação dos criminosos e ainda servir como mais um elemento probatório para a formação da convicção do Juízo. A Figura 1 apresenta o esquema gráfico da metodologia proposta.

Figura 1 - Metodologia proposta



Fonte: Elaborado pelos autores (2021)

Destaca-se que a primeira fase da metodologia se dá fora do processo de atribuição de autoria, pois é relativa à investigação policial. Dentro do domínio do processo proposto, a primeira ação de processamento é a análise estilométrica dos textos sem identificação do autor e sua caracterização quanto ao estilo de escrita. Cada texto com características singulares constituirá um subconjunto e a ele será atribuído um código identificador único, substituindo o nome do autor, até então desconhecido (Autor 1, Autor 2 etc.). Esses subconjuntos são armazenados (banco de dados) e serão usados para se comparar seus estilos de escrita com textos de autoria conhecida (Fases 3a e 3b). Na quarta fase, outro algoritmo faz a comparação entre as características estilométricas dos textos de autores conhecidos presentes na Internet (fóruns, redes sociais, blogs) e as características dos textos com conteúdo criminoso extraídos para a base de dados. Esse algoritmo pode também comparar as características estilométricas dos textos apreendidos em operações policiais com as características dos textos pedofílicos extraídos para a base de dados.

O *output* do processo (Fase 5) é a lista de suspeitos, ou seja, o processo proposto por esta metodologia mapeia as características estilométricas de textos suspeitos de serem criminosos e atribui a cada texto um código que ocupa a função do nome de autor, já que são desconhecidos os autores desse tipo de material. Sempre que o algoritmo encontra um texto com as mesmas características de escrita presentes em algum dos textos de autoria declarada, faz a correspondência, atribuindo ao texto que tinha a autoria inicialmente desconhecida o nome de um ou mais suspeitos.

6 CONCLUSÕES

Acredita-se que há uma aplicabilidade bastante viável da metodologia proposta na identificação de autores de material de incentivo ao abuso sexual de crianças e adolescentes via textos distribuídos pela Internet. Este estudo mostra que não há muitos trabalhos científicos sobre este assunto, o que sugere existir um campo ainda pouco explorado de pesquisas. O presente trabalho se apresenta como uma tentativa de ajudar a preencher essa lacuna, não só

apresentando uma sugestão de metodologia, mas também provocando a comunidade acadêmica a discutir as potencialidades da atribuição de autoria no combate à pedofilia e aos crimes a ela relacionados. A metodologia proposta foi pensada para auxiliar investigações de pedofilia pela Internet, mas pode ser aplicada para identificar outros tipos de criminosos, como *haters*, praticantes de bullying ou terroristas, o que se mostra como um de seus pontos fortes, além do ineditismo.

Contudo, há a necessidade de aprofundamento da pesquisa, a fim de se aprimorar as etapas e processos propostos. Além disso, a metodologia precisa passar por testes de adequação e validação, antes de estar disponível para a aplicação.

6.1 PERSPECTIVAS FUTURAS

A proposta deste trabalho necessita de uma expansão futura na qual algoritmos “treinados” sejam utilizados para identificar textos na Internet com autoria conhecida e que apresentem semelhança estilométrica com textos criminosos sem autoria conhecida. Construir e treinar um algoritmo assim é um grande desafio, porque essa abordagem contém em si todo o processo de análise estilográfica e toda a programação necessária para essa espécie de *Web Scraping* (raspagem de dados na Web). Além disso, há de se ter a preocupação com o treinamento desses algoritmos para que não retornem uma lista demasiado grande de suspeitos. Uma primeira medida para evitar isso é ter bem refinada a análise estilométrica dos textos contidos nos manuais de pedofilia apreendidos. Somado a isso, os algoritmos poderiam ser baseados em aprendizado de máquina, melhorando por experiência a análise estilométrica que fizerem, com o intuito de melhorar sua acuidade.

REFERÊNCIAS

ABBASI, A.; CHEN, H. Visualizing Authorship for Identification. In: Mehrotra, S., Zeng, D.D., Chen, H., Thuraisingham, B., Wang, FY. (eds) Intelligence and Security Informatics. ISI 2006. Lecture Notes in Computer Science, vol 3975. Springer, Berlin, Heidelberg. https://doi.org/10.1007/11760146_6.

ABBASI, A.; CHEN, H. Writeprints. **ACM Transactions on Information Systems**, [s. l.], v. 26, n. 2, p.1-29, 1 mar. 2008. DOI <http://dx.doi.org/10.1145/1344411.1344413>.

BHARGAVA, M.; MEHNDIRATTA, P.; ASAWA, K. Stylometric Analysis for Authorship Attribution on Twitter. **Big Data Analytics**, [s. l.], p. 37-47, 2013. Springer International Publishing. DOI http://dx.doi.org/10.1007/978-3-319-03689-2_3.

BRASIL. **Decreto-lei n.º 2.848, de 7 de dezembro de 1940**. Código penal. Disponível em: http://www.planalto.gov.br/ccivil_03/decreto-lei/del2848.htm. Acesso em: 19 maio 2021.

BRASIL. **Lei nº 8.069, de 13 de julho de 1990**. Dispõe sobre o Estatuto da Criança e do Adolescente e dá outras providências. Disponível em: http://www.planalto.gov.br/ccivil_03/leis/l8069.htm. Acesso em: 19 maio 2021.
CHILDHOOD BRASIL. **Quem somos**. 2021a. Disponível em: <https://www.childhood.org.br/quem-somos#intro>. Acesso em: 18 maio 2021.

CHILDHOOD BRASIL. **Números da causa**. 2021b. Disponível em: <https://www.childhood.org.br/nossa-causa#numeros-da-causa>. Acesso em: 18 maio 2021.

HADJIDJ, R.; DEBBABI, M.; LOUNIS, H.; IQBAL, F.; SZPORER, A.; BENREDJEM, D. Towards an integrated e-mail forensic analysis framework. **Digital Investigation**, v. 5, n. 3-4, p. 124-137, 2009. Disponível em: <https://www.sciencedirect.com/science/article/abs/pii/S1742287609000036>. Acesso em: 18 maio 2021.

ESCALANTE, H. J. *et al.* Early detection of deception and aggressiveness using profile-based representations. **Expert Systems with Applications**, v. 89, p. 99-111, 2017. DOI <https://doi.org/10.1016/j.eswa.2017.07.040>.

FRANCO, D. P.; MAGALHÃES, S. R. A dark web: navegando no lado obscuro da Internet. **Amazônia em Foco**, Castanhal, v. 4, n. 6, p. 18-33, jan./jul. 2015. Disponível em: <http://revista.fcat.edu.br/index.php/path/article/download/27/137>. Acesso em: 20 jan. 2019.

GE, Z.; SUN, Y.; SMITH, M. J. T. Authorship attribution using a neural network language model. **School of Electrical and Computer Engineering**, p. 4212–4213, 2016. Disponível em: <https://ojs.aaai.org/index.php/AAAI/article/view/9924/9783>. Acesso em: 20 jan. 2019.

ISHIHARA, S. A comparative study of likelihood ratio based forensic text comparison procedures: multivariate Kernel Density with Lexical Features vs. Word N-grams vs. Character N-grams. *In: CYBERCRIME AND*

TRUSTWORTHY COMPUTING CONFERENCE, 5., 2014. p. 1-11. Disponível em: <https://openresearch-repository.anu.edu.au/handle/1885/102627>. Acesso em: 16 jan. 2019.

ISHIHARA, S. Strength of linguistic text evidence: a fused forensic text comparison system. **Forensic Science International**, v. 278, p. 184-197, 2017. DOI 10.1016/j.forsciint.2017.06.040

KOURTIS, I.; STAMATATOS, E. Author identification using semi-supervised Learning Notebook for PAN at CLEF 2011. **University of the Aegean**, 2011.

MOREIRA, M. **Análise de manuais de pedofilia na dark web para prevenção de crimes sexuais contra crianças e adolescentes**. 2020. Dissertação (Mestrado em Ciência da Informação) – Programa de Pós-Graduação em Ciência da Informação, Universidade Federal de Santa Catarina, Florianópolis, SC, 2020.

PENG, F.; SCHUURMANS, D.; KESELJ, V.; WANG, S. Language independent authorship attribution using character level language models. *In*: CONFERENCE ON EUROPEAN CHAPTER OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS, 10., (EACL '03), v. 1, 2003. **Proceedings** [...]. Association for Computational Linguistics, USA, 2003. p. 267–274.

RAMNIAL, H.; PANCHOO, S.; PUDARUTH, S. Authorship attribution using stylometry and machine learning techniques. **Advances in Intelligent Systems And Computing**, [s. l.], p.113-125, 29 ago. 2015. DOI http://dx.doi.org/10.1007/978-3-319-23036-8_10.

ROCHA, A.; SCHEIRER, W.; FORSTALL, C.; CAVALCANTE, T.; THEOPHILO, A.; SHEN, B.; CARVALHO, A.; STAMATATOS, E. Authorship Attribution for Social Media Forensics. **IEEE Transactions on Information Forensics and Security**, [s. l.], v. 12, n. 1, p.121-122, jan. 2017. DOI 10.1109/TIFS.2016.2603960.

SAFARNET. **Institucional**. Disponível em: <https://new.safarnet.org.br/content/institucional#mobile>. Acesso em: 18 maio 2021.

STAMATATOS, E. A survey of modern authorship attribution methods. **Journal of the American Society for Information Science and Technology**, [s. l.], v. 60, n. 3, p.538-556, mar. 2009. DOI <http://dx.doi.org/10.1002/asi.21001>.

VILLAR-RODRIGUEZ, E. *et al.* A feature selection method for author identification in interactive communications based on supervised learning and language typicality. **Engineering Applications of Artificial Intelligence**, v. 56, p. 175-184, 2016.

YANG, M.; CHOW, K. Authorship Attribution for Forensic Investigation with

Thousands of Authors. **ICT Systems Security and Privacy Protection**, [s. l.], p.339-350, 2014. DOI http://dx.doi.org/10.1007/978-3-642-55415-5_28.

ZHENG, R.; LI, J.; CHEN, H.; HUANG, Z. A framework for authorship identification of online messages: Writing-style features and classification techniques. **Journal of the American Society for Information Science and Technology**, [s. l.], v. 57, n. 3, p. 378-393, 2006. DOI <http://dx.doi.org/10.1002/asi.20316>

XYLOGIANNOPOULOS, K.; KARAMPELAS, P.; ALHAJJ, R. Text mining for plagiarism detection: multivariate pattern detection for recognition of text similarities. *In: IEEE/ACM INTERNATIONAL CONFERENCE ON ADVANCES IN SOCIAL NETWORKS ANALYSIS AND MINING (ASONAM)*, 2018., Barcelona, Spain, p. 938-945, ago. 2018. DOI <http://dx.doi.org/10.1109/asonam.2018.8508265>.

AN AUTHORSHIP ATTRIBUTION MODEL APPLIED TO PEDOPHILIA CRIME INVESTIGATIONS

ABSTRACT

Objectives: Identify the current state of the art of scientific research in the field of authorship attribution applied to investigations of sexual crimes against children and adolescents over the Internet involving written material. Propose a methodology for using authorship attribution to identify suspected authors of texts with content that encourages child and adolescent sexual abuse. **Methodology:** This is a qualitative research that uses the Systematic Review of Literature to identify works that deal with the techniques of authorship attribution in order to seek scientific evidence of its application to problems similar to the one addressed in the present study. **Results:** The current state of the art of scientific research that relates the use of authorship attribution techniques to texts on the internet that encourage the practice of sexual abuse of children and adolescents is presented and, from this, a methodology is proposed to identification of authors of texts with those characteristics. **Conclusions:** It is concluded that there is not an abundance of scientific research on this topic, which suggests that it is an open field for further studies. It is also concluded that it is fully possible to apply the techniques of authorship attribution in the identification of the probable authors of texts that aim to guide and encourage the practice of child and adolescent sexual abuse, which is explained by the proposed methodology.

Descriptors: Child and adolescent sexual abuse on the internet. Authorship attribution. Stylogometry. Pedophilia. Police investigation.

UNA METODOLOGÍA DE ASIGNACIÓN DE AUTORIZACIÓN APLICADA A INVESTIGACIONES SOBRE ABUSO SEXUAL INFANTIL

RESUMEN

Objetivos: Identificar el estado actual de la investigación científica en el campo de la atribución de autoría aplicada a las investigaciones de delitos sexuales contra niños, niñas y adolescentes a través de Internet que involucran material escrito. Proponer una metodología para el uso de la atribución de autoría para identificar a los presuntos autores de textos con contenido que fomenta el abuso sexual de niños y adolescentes.

Metodología: Se trata de una investigación cualitativa que utiliza la Revisión Sistemática de Literatura para identificar trabajos que traten sobre las técnicas de atribución de autoría con el fin de buscar evidencia científica de su aplicación a problemas similares al abordado en el presente estudio. **Resultados:** se presenta el estado actual de la investigación científica que relaciona el uso de técnicas de atribución de autoría a textos presentes en internet que incentivan la práctica del abuso sexual de niños, niñas y adolescentes y, a partir de ello, se propone una metodología para la identificación de autores de textos con esas características. **Conclusiones:** Se concluye que no existe una abundancia de investigaciones científicas sobre este tema, lo que sugiere que es un campo abierto para estudios posteriores. También se concluye que es plenamente posible aplicar las técnicas de atribución de autoría en la identificación de los probables autores de textos que pretenden orientar y fomentar la práctica del abuso sexual infantil y adolescente, lo cual fue explicado por la metodología propuesta.

Descriptores: Abuso sexual infantil en Internet. Atribución de autoría. Estilometría. Pedofilia. Investigación policial.

Recibido em: 21.05.2021

Aceito em: 30.04.2022