

# Minerando e Caracterizando Dados de Currículos Lattes

Luciano A. Digiampietri<sup>1</sup>, Jesús P. Mena-Chalco<sup>2</sup>, José J. Pérez-Alcázar<sup>1</sup>,  
Esteban F. Tuesta<sup>1</sup>, Karina V. Delgado<sup>1</sup>, Rogério Mugnaini<sup>1</sup>, Gabriela S. Silva<sup>1</sup>

<sup>1</sup>Escola de Artes, Ciências e Humanidades da Universidade de São Paulo

<sup>2</sup>Centro de Matemática, Computação e Cognição da Universidade Federal do ABC

digiampietri@usp.br, jesus.mena@ufabc.edu.br

**Abstract.** *Curricula from the Lattes Platform are a vast source of information for the creation and analyzes of researcher's social networks. Due to the great amount of data, the manual filling, and the use of semi-structured data there are several challenges in the use of these curricula. This paper presents a database produced from the mining of more than one million Lattes curricula. Moreover, it describes some characteristics and relations of these curricula, directions and challenges to the production and analyzes of social networks generated from these data.*

**Resumo.** *Currículos da Plataforma Lattes são uma vasta fonte de informação para a criação e análise de redes sociais de pesquisadores. Devido à quantidade de dados, ao preenchimento manual e ao uso de dados semiestruturados existem diversos desafios para a utilização destes currículos. Este artigo apresenta um banco de dados produzido a partir da mineração de mais de um milhão de Currículos Lattes, descrevendo algumas características e relações desses currículos, direções e desafios para a produção e análise de redes sociais a partir destes dados.*

## 1. Introdução

Atualmente, é possível encontrar na Web uma grande quantidade de dados referentes aos mais diversos assuntos. Dentre esses dados estão informações muito relevantes aos pesquisadores, como publicações científicas, informações sobre projetos de pesquisa e mesmo currículos de pesquisadores.

Ao se tratar de dados referentes à pesquisa, o Brasil apresenta uma característica peculiar: a existência de um cadastro nacional de currículos de pesquisadores, o Currículo Lattes, que congrega informações sobre publicações, orientações, projetos de pesquisa, entre outras. O Currículo Lattes foi lançado e padronizado em agosto de 1999 pelo CNPq como sendo o formulário de currículo a ser utilizado no âmbito do Ministério da Ciência e Tecnologia e CNPq<sup>1</sup>, e no ano de 2007 ultrapassou a marca de 1 milhão de currículos.

Currículos da Plataforma Lattes são uma vasta fonte de informação para a criação e análise de redes sociais de pesquisadores [Balancieri et al. 2005]. Devido à quantidade de dados, ao preenchimento manual e ao uso de dados semiestruturados existem diversos desafios computacionais para a utilização destes currículos. Esse grande volume de

---

<sup>1</sup><http://lattes.cnpq.br/conteudo/historico.htm>

informações tem sido pouco usado, servindo, tipicamente, para avaliar (ou verificar) dados de pesquisadores individualmente ou de pequenos grupos de pesquisadores.

Uma busca sobre “análise de currículo(s)” (“*analysis of curric\**”) na Web of Science recupera 43 referências a artigos científicos, sendo 50% deles na área de educação. Neste conjunto, evidencia-se o interesse em estudos sobre competência profissional e formação acadêmica. Com a Plataforma Lattes, abre-se a oportunidade de se investigar questões de interesse da área de Política Científica, possibilitando-se estudar não apenas a produtividade dos pesquisadores, mas suas relações de colaboração, vinculando assim dois temas normalmente abordados separadamente: Análise de Redes Sociais e Cientometria. Porém, o acesso à base Lattes ainda é pouco facilitado pelo CNPq, fazendo com que estudos em nível macro, ou utilizem a interface do Censo do Diretório dos Grupos de Pesquisa no Brasil [Guimarães 2004] (via que não permite muita liberdade para análise dos dados), ou dependam de solicitação dos dados ao CNPq [Leite et al. 2011, Mugnaini et al. 2011] (acesso que nem sempre é garantido).

Este artigo visa a apresentar um banco de dados formado por mais de um milhão de currículos minerados da Plataforma Lattes e que foram processados, organizados e analisados para servirem de base para a produção e análise de redes sociais de pesquisa. Neste artigo, diversas características desse banco serão apresentadas, bem como algumas relações identificadas entre os currículos, as quais ligam orientandos e orientadores e coautores de publicações científicas. Além disso, são discutidos alguns cuidados que devem ser considerados nos trabalhos que analisam currículos da Plataforma Lattes.

O restante deste artigo está organizado da seguinte forma. A Seção 2 apresenta os trabalhos correlatos. A Seção 3 descreve os processos de mineração e produção do banco de dados. A Seção 4 contém a caracterização, análise e cuidados relacionados ao banco de dados. Por fim, a Seção 5 contém as considerações finais e os trabalhos futuros.

## 2. Trabalhos Correlatos

Há mais de 60 anos existem estudos sobre análise de redes sociais, havendo uma vasta literatura sobre o assunto [Newman 2001]. Esta seção apresenta apenas uma visão geral dos trabalhos correlatos com enfoque na análise de redes sociais de pesquisadores e, especialmente, naqueles relacionados ao uso de Currículos Lattes.

Silva e Smit [Silva and Smit 2009] avaliaram a organização e qualidade da informação científica disponível nos Currículos Lattes. Eles concluíram que há bastante comprometimento no preenchimento dos currículos, por mais que existam diversos pequenos erros de preenchimento. Por analisar uma amostra não muito grande de currículos, eles não calcularam estatísticas sobre a Plataforma Lattes.

Mena Chalco e César Júnior [Mena-Chalco and Cesar Junior 2009] desenvolveram e disponibilizaram uma ferramenta chamada *scriptLattes*<sup>2</sup>. Esta ferramenta recebe uma lista de identificadores de Currículos Lattes e gera diversas páginas HTML, organizando as informações dos currículos (tanto sumarizando informações quanto separando as informações por categoria). Além disso, a ferramenta gera um mapa da distribuição geográfica dos pesquisadores envolvidos na análise e um grafo de coautorias.

---

<sup>2</sup><http://scriptlattes.sourceforge.net/>

Alves et al. [Alves et al. 2011] desenvolveram um sistema chamado SUCUPIRA, uma ferramenta online que além de possuir funcionalidades específicas para baixar e organizar currículos Lattes une algumas ferramentas/APIs para visualizar informações de currículos e suas redes de coautoria. O sistema foi testado analisando-se o conjunto de autores de um programa de pós-graduação da UFMG.

Digiampietri e Silva [Digiampietri and da Silva 2011] desenvolveram uma infraestrutura para encontrar currículos de grupos de pesquisadores. Dada uma lista com o nome dos pesquisadores do grupo de interesse, a infraestrutura procura pelos currículos utilizando as APIs de busca da Google<sup>3</sup> e da Microsoft<sup>4</sup>, baixa esses currículos, gera as redes de coautoria e analisa a produção dos pesquisadores cruzando as informações de cada currículo com as informações dos documentos de área da CAPES.

O banco de dados apresentado neste artigo se diferencia dos demais trabalhos por apresentar um vasto conjunto de currículos (ao nosso entender, muito maior do que os analisados nos demais trabalhos). Outra característica relevante é a granularidade da informação dos currículos, pois neste banco as informações são organizadas campo a campo (e não publicação a publicação, por exemplo). Além disso, os dados foram cuidadosamente estruturados e inseridos em um banco de dados relacional visando a facilitar consultas, análises e enriquecimento dos dados. Por se tratar de conjunto de dados representativo, foi possível analisar algumas características e gerar algumas estatísticas sobre os Currículos Lattes, como será apresentado nas próximas seções.

### 3. Banco de Dados de Currículos

Esta seção descreve a busca, obtenção e organização dos dados obtidos de currículos da Plataforma Lattes.

#### 3.1. Busca pelos Currículos

Os currículos da Plataforma Lattes foram obtidos pela internet, utilizando-se o comando *wget* para baixar cada um dos currículos. Para isto, foi necessário encontrar o identificador numérico de cada currículo, pois este é necessário para compor a URL (*Uniform Resource Locator*) completa do currículo, a qual é parâmetro da ferramenta *wget*.

Para encontrar os identificadores dos currículos duas estratégias foram usadas. Na primeira, foram feitas consultas na interface de busca provida pela Plataforma Lattes. Mais especificamente, foram feitas 80 consultas, cada uma utilizando como palavras-chave as (sub)áreas de conhecimento da própria plataforma. Foi desenvolvido um *parser* para encontrar os identificadores dos currículos de cada uma das respostas às consultas realizadas. Esta estratégia permitiu a identificação de centenas de milhares de currículos. Estes currículos foram baixados e examinados automaticamente em busca de identificadores de outros currículos (por exemplo, de coautores, orientadores ou orientandos). Este exame constituiu a segunda estratégia de busca. Sempre que um novo currículo era identificado, este era baixado e examinado da mesma forma que os demais.

Combinando-se estas duas estratégias 1.236.548 currículos foram baixados totalizando pouco mais de 16 GB em arquivos HTML. É importante destacar que as estratégias

---

<sup>3</sup>[www.google.com](http://www.google.com)

<sup>4</sup>[www.bing.com](http://www.bing.com)

adotadas não visavam à obtenção de toda a base de currículos da Plataforma Lattes, mas sim um conjunto significativo de currículos para serem processados e servirem de base para a criação e análise de redes sociais de pesquisadores. Levando-se em conta que o CNPq anunciou que no ano de 2007 a base de currículos atingiu um milhão de currículos, considera-se que foi obtida uma quantidade representativa de currículos.

Todos os currículos foram baixados em maio de 2011, a partir desta data foram feitos apenas processamentos com base nas informações dos currículos já baixados.

### 3.2. Processamento Inicial dos Currículos

Cada um dos currículos baixados foi submetido a duas etapas iniciais de processamento. Na primeira, foram retirados todos os caracteres especiais, espaços em brancos excedentes e fins de linha, de forma a facilitar a identificação de informações executada na etapa subsequente de processamento.

Na segunda etapa de processamento, foram desenvolvidas expressões regulares para dividir cada currículo em suas seções principais (Dados Gerais; Linhas de pesquisa; Projetos; Áreas; Produção em C, T & A; Bancas; Eventos; e Orientações). As informações de cada uma destas seções foram processadas por expressões regulares específicas para a identificação dos itens e campos de interesse. Por exemplo, informações referentes a Produção em C, T & A foram subdivididas em Produção bibliográfica; Produção técnica; Produção artística/cultural; e Demais trabalhos. Informações sobre a Produção bibliográfica foram subdivididas em 7 categorias: artigos completos publicados em periódicos; artigos aceitos para publicação; trabalhos completos publicados em anais de congressos; resumos expandidos publicados em anais de congressos; resumos publicados em anais de congressos; livros publicados organizados ou edições; e capítulos de livros publicados. Expressões regulares foram desenvolvidas para identificar os campos de cada uma destas categorias. Apenas para exemplificar, os seguintes campos foram extraídos de cada artigo completo publicado em periódico: título, local, páginas, volume, autores, ano de publicação, nome do periódico, número e ISSN. Para cada currículo processado foi produzido um arquivo XML para estruturar as informações identificadas.

### 3.3. Banco de Dados Propriamente Dito

Para organizar as informações obtidas dos Currículos Lattes, um banco de dados foi criado dentro do SGBD PostgreSQL. O esquema deste banco de dados pode ser observado na Figura 1. Cada uma das tabelas será descrita de maneira sucinta a seguir. Neste artigo, cada pessoa que possui um currículo Lattes será chamada de pesquisador.

**Curriculos:** tabela contendo as informações gerais de cada currículo: identificador (*lattesID*), nome, tipo de bolsa produtividade, sexo e última atualização do currículo.

**NomesUsadosEmCitacoes:** tabela com os nomes/abreviações usadas pelo autor em suas citações.

**Formacoes:** tabela contendo dados sobre a formação de uma dada pessoa, incluindo o período de formação, o título, a instituição, o identificador do orientador e o nome do orientador.

**AreasDeAtuacao:** contém a lista das áreas de atuação informadas pelo pesquisador (contendo grande área, área, subárea e especialidade).

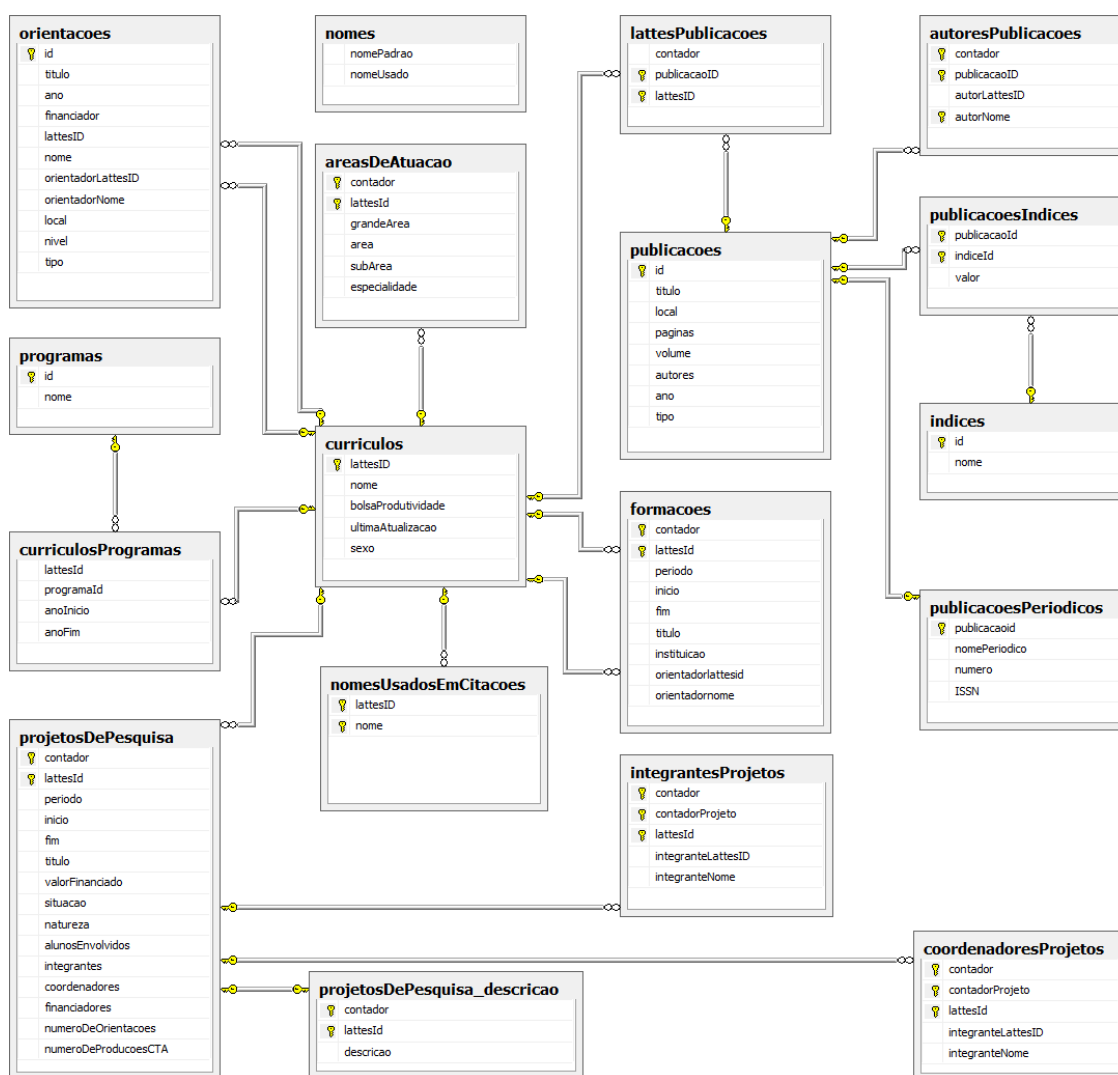
**ProjetosDePesquisa:** lista dos projetos de pesquisa no qual cada pesquisador está envolvido.

**ProjetosDePesquisa\_descricao:** descrição dos projetos de pesquisa, na forma de texto livre.

**CoordenadoresProjetos:** tabela para listar os coordenadores de cada projeto de pesquisa.

**IntegrantesProjetos:** tabela para listar os integrantes de cada projeto de pesquisa.

**Publicacoes:** tabela para armazenar os registros de todas as publicações de todos os currículos. Há sete tipos de publicações que estão sendo consideradas: artigos completos publicados em periódicos; artigos aceitos



**Figura 1. Diagrama Entidade Relacionamento do Banco de Dados Desenvolvido**

para publicação; trabalhos completos publicados em anais de congressos; resumos expandidos publicados em anais de congressos; resumos publicados em anais de congressos; livros publicados organizados ou edições; e capítulos de livros publicados.

**LattesPublicacoes:** tabela para vincular cada currículo ao conjunto de publicações cadastradas pelo pesquisador.

**AutoresPublicacoes:** tabela com o conjunto de (co)autores de cada publicação.

**Orientacoes:** lista das orientações feitas por cada pesquisador. Há sete tipos de orientação que estão sendo consideradas: Orientações de Pós-Doutorado; Teses de Doutorado; Orientações de Outra Natureza; Dissertações de Mestrado; Monografias de Conclusão de Curso de Aperfeiçoamento; Iniciações Científicas; e Trabalhos de Conclusão de Graduação.

Todos os arquivos XML produzidos foram lidos e seus conteúdos inseridos no banco de dados. Um total de 88.402.026 de registros foram cadastrados. A Tabela 1 apresenta o total de registros cadastrados em cada uma das tabelas supracitadas.

Além dessas tabelas, algumas tabelas adicionais foram criadas para permitir ou facilitar análises mais sofisticadas sobre o conjunto de dados. Por exemplo, a tabela *Nomes* contém os nomes na forma que foram usados pelo pesquisador e uma versão canônica

**Tabela 1. Número de Registros em Cada uma das Tabelas do Banco de Dados**

<b>Tabela</b>	<b>Número de Registros</b>
projetosdepesquisa_descricao	1.069.884
curriculos	1.236.548
nomesusadosemcitacoes	1.353.467
projetosdepesquisa	1.378.885
coordenadoresprojetos	1.380.087
formacoes	3.250.846
areasdeatuacao	3.256.019
orientacoes	4.329.993
integrantesprojetos	4.915.223
publicacoes	11.529.218
lattespublicacoes	11.529.218
autorespublicacoes	43.172.638
<b>Total</b>	<b>88.402.026</b>

(sem acentos e caracteres especiais) para facilitar a comparação entre nomes registrados pelo pesquisador e por seus coautores, orientandos ou orientadores. As tabelas *Indices* e *PublicacoesIndices* foram criadas para armazenar informações adicionais sobre cada publicação, por exemplo, JCR e número de citações (informações que não necessariamente estão presentes nos currículos Lattes, mas que serão obtidas numa etapa futura de enriquecimento da base de dados). Além disso, um campo adicional foi inserido na tabela *Publicacoes* chamado *idUnico* que está sendo usado para relacionar diferentes registros de publicações que se referem à mesma entidade (por exemplo, três coautores registraram o mesmo artigo, então há três registros deste artigo na base e o campo *idUnico* conterá o mesmo valor para esses registros). Esta informação também não pode ser obtida diretamente/explicitamente dos currículos mas, como será visto na Subseção 4.2, ela poderá ser inferida das informações presentes nos currículos.

#### **4. Caracterização do Banco de Dados de Currículos**

Esta seção contém a caracterização da distribuição dos dados nos currículos, uma breve análise sobre dados e um conjunto de dicas e cuidados que devem ser considerados na análise de Currículos Lattes.

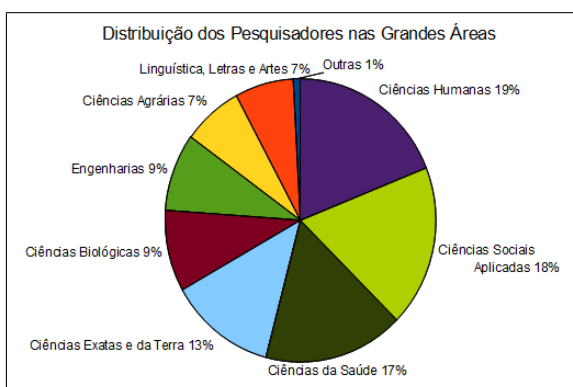
##### **4.1. Descrição do Banco de Dados**

Nesta seção serão apresentadas algumas características do banco de dados formado com as informações dos currículos analisados.

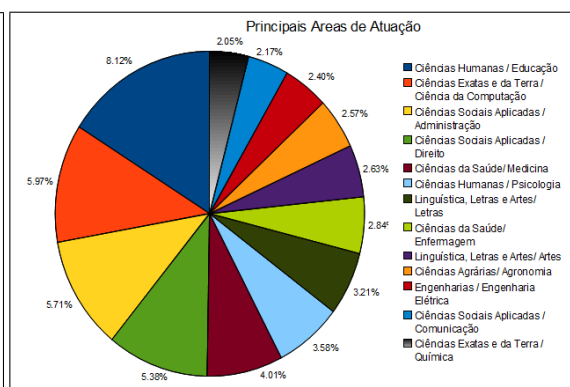
Na média, cada pesquisador informou que atua em 2,61 grandes áreas, áreas ou subáreas diferentes. Em mais detalhes, 263.775 pesquisadores informaram que atuam em mais de uma grande área (das oito grandes áreas disponíveis para seleção na Plataforma Lattes). A Figura 2 apresenta a distribuição dos pesquisadores nas grandes áreas.

As grandes áreas são divididas em dezenas de áreas, porém as 13 áreas mais informadas pelos pesquisadores correspondem, juntas, a mais de 50% do total de declarações de áreas de atuação. A Figura 3 apresenta a distribuição dos pesquisadores nessas áreas.

Conforme apresentado, todos os currículos foram baixados em maio de 2011. Enquanto uma grande parcela dos currículos foi atualizada nos últimos doze meses (contados

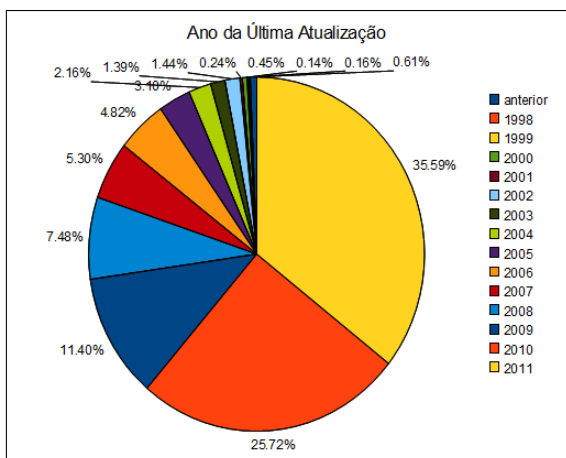


**Figura 2. Distribuição nas Grandes Áreas de Atuação**

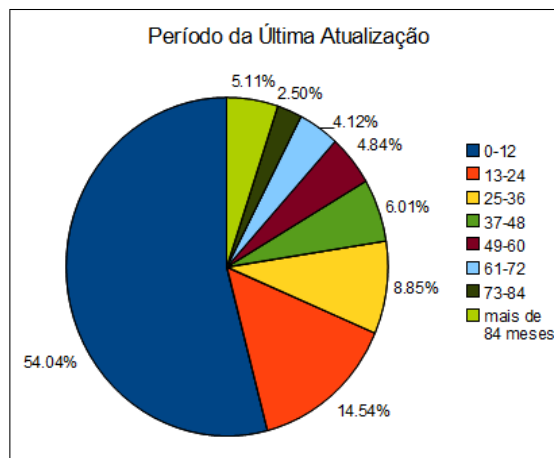


**Figura 3. Principais Áreas de Atuação**

de junho de 2010 a maio de 2011), há diversos currículos que estão desatualizados há diversos anos. Na média, cada currículo foi atualizado pela última vez há 27 meses, mas a mediana é de apenas 10 meses. As Figuras 4 e 5 apresentam, respectivamente, o ano da última atualização dos currículos e o período (em intervalos de doze meses contados a partir de maio de 2011) da última atualização. Pode-se observar que mais de 75% dos currículos foram atualizados nos últimos 36 meses.



**Figura 4. Ano da Última Atualização do Currículo**

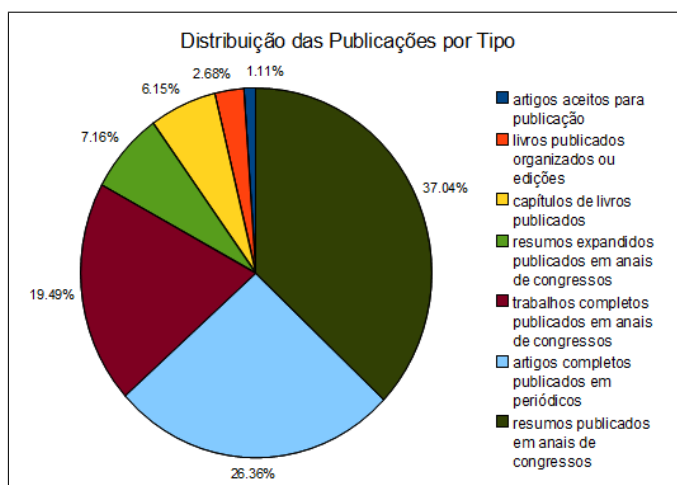


**Figura 5. Período da Última Atualização do Currículo**

Nos currículos analisados foram encontrados 11.529.218 de registros de publicações, ou seja, uma média de 9,32 publicações por currículo. É importante lembrar que esses mais de 11,5 milhões de registros de publicações possuem redundâncias, pois diferentes coautores inserem uma mesma publicação em seus currículos.

As publicações estão organizadas em sete tipos, conforme pode ser observado na Figura 6. Na média, cada publicação tem 3,74 autores. A Figura 7 apresenta as médias de autores para cada um dos tipos de publicação.

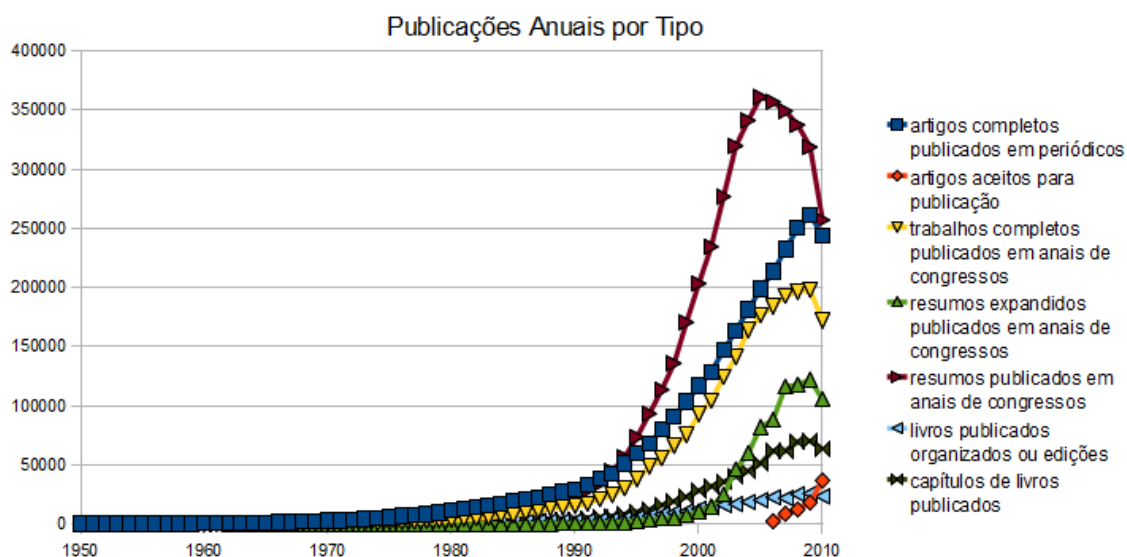
A Figura 8 apresenta a evolução do número de registros de publicações cadastrados nos currículos ao longo dos anos. Nesta figura podemos observar um crescimento no número de publicações ano a ano. Nos últimos três anos esse crescimento não pode ser



**Figura 6. Distribuição das Publicações por Tipo**



**Figura 7. Número Médio de Autores**



**Figura 8. Evolução no Número de Publicações no Tempo**

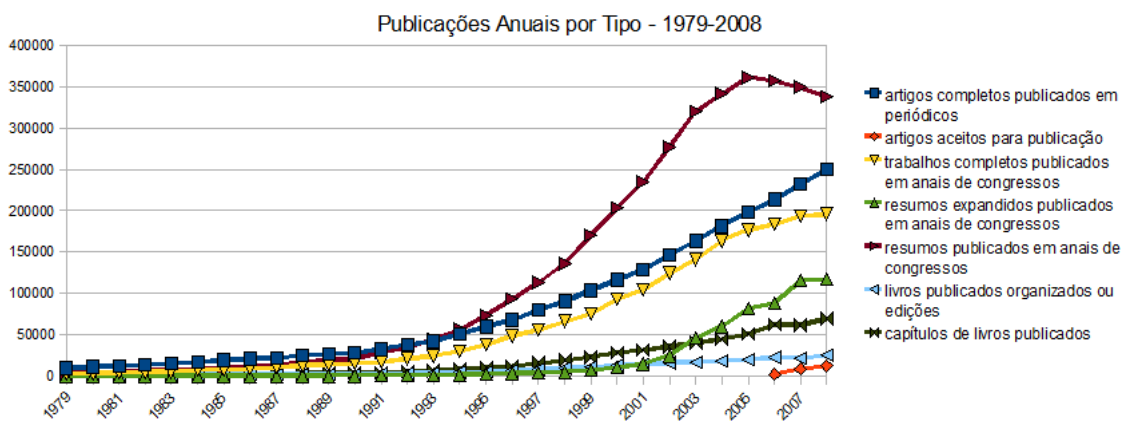
observado, mas isto ocorre devido à não atualização dos currículos (que na média foram atualizados pela última vez há 27 meses).

Uma visualização melhor no número de registros de publicações pode ser feita considerando-se apenas os trinta anos de 1979 até 2008, período no qual a grande maioria dos currículos está atualizada. A Figura 9 apresenta estas informações. Ao observarmos os registros de artigos publicados em periódicos é possível constatar um crescimento exponencial no número desses registros. De fato, nessa janela de trinta anos, o número desses registros cresceu cerca de 12% ao ano.

Dos 1.236.548 currículos cadastrados, 13.797 possuem Bolsa Produtividade do CNPq (1,12% dos pesquisadores). A Figura 10 apresenta a distribuição das bolsas em seus sete níveis.

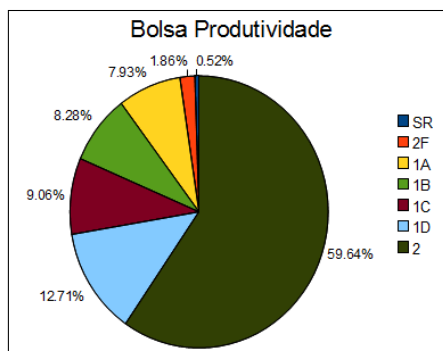
No banco de dados produzido, há 4.329.993 orientações cadastradas (cerca de 3,5



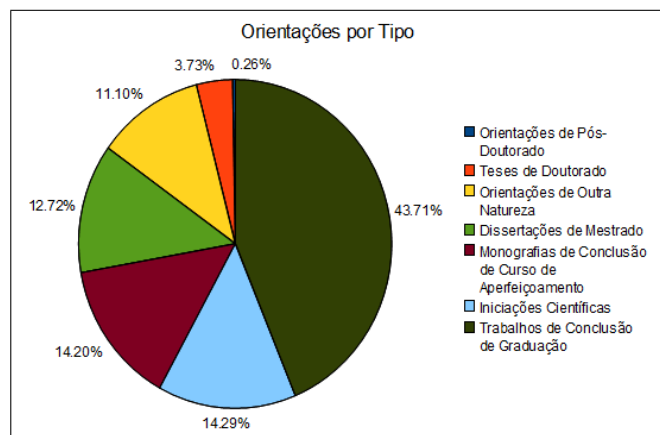


**Figura 9. Evolução das Publicações no Período de 1979 à 2008**

orientações por currículo). Estas orientações estão divididas em sete categorias, conforme ilustrado na Figura 11. As orientações de mestrado e doutorado correspondem a pouco menos de 16,5% do total de orientações.



**Figura 10. Distribuição das Bolsas Produtividade**

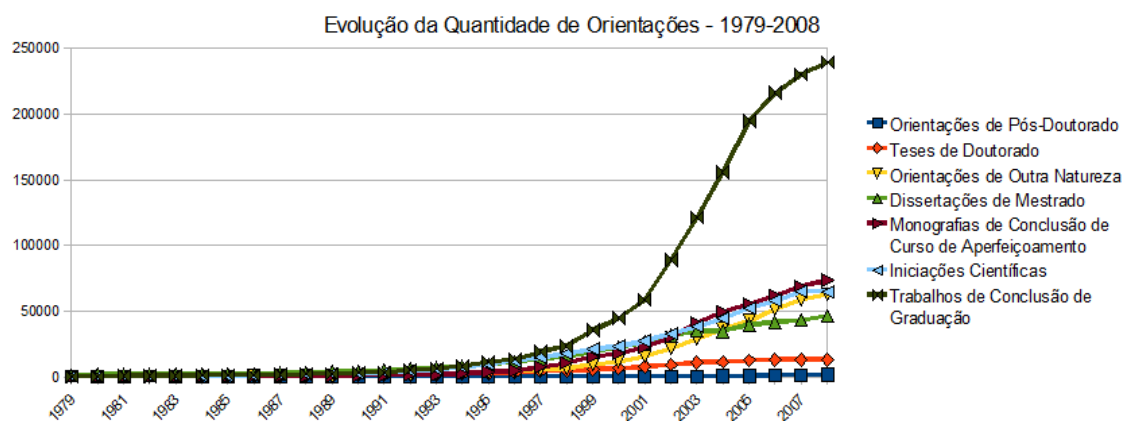


**Figura 11. Distribuição das Orientações por Categoria**

A Figura 12 apresenta a evolução da quantidade de orientações de cada tipo no período de trinta anos, de 1979 a 2008. Duas informações merecem destaque nesta evolução: a primeira é a grande quantidade de trabalhos de conclusão de curso que vêm sendo cadastrados nos últimos anos. Conforme apresentado na Figura 11, este tipo de orientação corresponde a mais de 43% do total de orientações. A segunda informação é a relação entre o número de orientações de mestrado e doutorado. Considerando esta janela de trinta anos, nos dez primeiros anos de análise, para cada 5,4 orientações de mestrado havia uma orientação de doutorado. Ao considerar os dez anos seguintes essa relação caiu para 3,9 para 1. Observando os dez últimos anos desse período (de 1999 a 2008) a relação caiu para menos de 3,3 para 1, indicando um crescimento proporcional bem maior no número de orientações de doutorado em relação às orientações de mestrado.

#### 4.2. Relações entre Currículos

Uma das grandes motivações da criação do banco de dados de Currículos Lattes é a geração e análise das redes sociais de pesquisadores. Para a montagem dessas redes é



**Figura 12. Quantidade de Orientações por Tipo ao Longo dos Anos**

necessário estabelecer quais características serão consideradas para relacionar diferentes currículos.

Nesta seção serão apresentadas algumas das características mais óbvias para a criação dessas redes, bem como a quantidade de relações que já foram obtidas considerando cada uma delas.

A primeira relação obtida é a de coautorias, ou seja, dois ou mais pesquisadores que são coautores de uma mesma publicação. Na base de dados há 11.529.218 registros de publicações, porém vários desses registros representam a mesma publicação referenciada em diferentes currículos. Uma consulta simples à base de dados, procurando apenas por publicações com o mesmo título e do mesmo tipo permite identificar 1.843.464 de relações de coautoria. Este critério de publicações com o mesmo título e do mesmo tipo não é robusto o suficiente para garantir que dois registros diferentes se referam à mesma publicação, bem como haverá registros referentes à mesma publicação com títulos diferentes (devido a erros no preenchimento, excesso ou falta de espaços ou pontuações, etc.). Assim, é necessário o desenvolvimento de métodos eficientes e eficazes de resolução de entidades. De qualquer forma, esse valor de mais de 1,8 milhão de relações identificadas é bastante relevante.

Uma segunda maneira de utilizar a relação de coautorias é através da verificação dos nomes e identificadores dos autores de cada uma das publicações. Essa relação, da maneira que está disponível nos Currículos Lattes não liga dois registros de publicações, mas relaciona a publicação de um currículo com outro currículo, servindo de base para a ligação entre currículos. Os 11.529.218 de registros de publicações possuem, ao todo, 43.172.638 registros de (co)autores. Destes, há identificadores de currículos de 21.102.745 (co)autores, sendo que 11.529.218 são dos donos dos currículos nos quais a publicação foi informada. Assim, restam 9.573.527 relações de coautoria identificadas.

Outra relação relevante quando analisados os currículos de pesquisadores é a relação orientador/orientando. Ao se cruzar informações de formação de um pesquisador com as informações de orientação de outro foi possível identificar um total de 210.112 relações de orientação de doutorado. Aqui, só foram consideradas as relações onde foi possível fazer a verificação dupla: um pesquisador informou que foi orientado por um segundo pesquisador e este segundo informou que foi orientador do primeiro.

É possível também relacionar os pesquisadores por áreas (grandes áreas, áreas, subáreas ou especialidades) de interesse. Devido ao fato de, na média, cada pesquisador informar cerca de 2,6 áreas de interesse, só esta característica originaria milhões de relações entre pesquisadores.

### **4.3. Cuidados e Problemas Relacionados aos Dados**

Os currículos preenchidos na Plataforma Lattes possuem certa padronização imposta pelos formulários da plataforma e algumas verificações de valores (por exemplo, o ano de conclusão de um projeto não pode ser anterior ao ano de início).

Porém, algumas das verificações só foram introduzidas nos últimos anos, o que permite que informações mais antigas estejam inconsistentes. Além disso, grande parte da informação é preenchida manualmente, o que possibilita a ocorrência de diversos tipos de problemas que precisam ser tratados durante o processamento e análise dos currículos.

Um primeiro cuidado que deve ser considerado é a existência de currículos de homônimos na Plataforma Lattes, isto por si só não é um problema, mas precisa ser considerado sempre que o nome do pesquisador for utilizado para o estabelecimento de qualquer tipo de relação entre currículos. No banco de dados produzido neste artigo foram encontrados 22.169 currículos de homônimos (1,79% dos currículos do banco).

Um problema que ocorria principalmente nas versões mais antigas da Plataforma Lattes é o preenchimento incompleto das informações. Em suas primeiras versões havia uma menor quantidade de campos obrigatórios e de mecanismos de verificação, assim não é incomum encontrar um registro de publicação sem o nome de nenhum autor; publicações em revista sem nenhuma indicação do nome da revista, e assim por diante.

Um problema comum e cuja solução é difícil de ser totalmente automatizada é o preenchimento incorreto de informações. Por exemplo, é comum que coautores preencham os campos referentes a uma mesma publicação de maneiras diferentes, variando a grafia do título do artigo, veículo de publicação, nome dos coautores e assim por diante. Além disso, diferentes autores classificam a mesma publicação de diferentes maneiras (por exemplo, um autor diz que uma publicação é um resumo e outro diz que é um resumo expandido). Desta forma, o processamento das informações dos currículos deve incluir mecanismos de casamento aproximado de *strings* e/ou análise de padrões.

Outra característica dos dados de Currículos Lattes é a diferença entre as atualizações dos currículos. Alguns autores atualizam seus currículos mensalmente enquanto outros atualizam menos de uma vez por ano. Mesmo um currículo atualizado recentemente pode conter dados desatualizados. Por exemplo, o banco de dados contém mais de 10.000 artigos classificados como *aceitos para publicação* entre 2006 e 2007, sendo que a maioria dos currículos foi atualizada depois de maio de 2010.

## **5. Considerações Finais e Trabalhos Futuros**

Os currículos da Plataforma Lattes contêm uma quantidade muito grande e diversificada de informações que podem ser utilizadas como base para a construção e análise de redes sociais de pesquisa, sendo uma das bases de pesquisadores mais completa do mundo.

Enquanto algumas relações podem ser extraídas diretamente desta plataforma, há diversas outras que podem ser obtidas através de algoritmos de resolução de enti-

dades. O banco de dados também pode ser enriquecido com informações relacionadas a publicações (número de citações de cada artigo, por exemplo); dos veículos de publicação (obtenção de índices como JCR e SJR); e dados derivados dessas informações como índices G e H.

Além disso, cada pesquisador pode ser caracterizado por diferentes índices como *Page Rank* ou *Author Rank*, ou novos índices podem ser criados exclusivamente para analisar pesquisadores através de sua produção científica e/ou orientações.

Neste artigo foram apresentados os primeiros passos na direção de uma análise ampla dos dados de Currículos Lattes, das relações entre esses dados e das redes formadas pelos seus pesquisadores. Como trabalhos futuros, pretende-se desenvolver algoritmos robustos para a resolução de entidades de forma a determinar uma maior quantidade de relações de orientação e coautoria. Pretende-se também criar e analisar redes de pesquisadores considerando diferentes relações e diferentes métricas de redes.

### **Agradecimentos**

O trabalho apresentado neste artigo foi parcialmente financiado pelo Programa de Educação Tutorial (MEC/SESu).

### **Referências**

- Alves, A., Yanasse, H., and Soma, N. (2011). Sucupira: A system for information extraction of the lattes platform to identify academic social networks. In *Information Systems and Technologies (CISTI), 2011 6th Iberian Conference on*, pages 1–6.
- Balancieri, R., Bovo, A. B., Kern, V. M., Pacheco, R. C. d. S., and Barcia, R. M. (2005). AA análise de redes de colaboração científica sob as novas tecnologias de informação e comunicação: um estudo na Plataforma Lattes. *Ciência da Informação*, 34:64 – 77.
- Digiampietri, L. A. and da Silva, E. E. (2011). A framework for social network of researchers analysis. *Iberoamerican Journal of Applied Computing*, 1(1):1 – 24.
- Guimarães, J. A. (2004). A pesquisa médica e biomédica no Brasil. Comparações com o desempenho científico brasileiro e mundial. *Ciência & Saúde Coletiva*, 9:303 – 327.
- Leite, P., Mugnaini, R., and Leta, J. (2011). A new indicator for international visibility: exploring brazilian scientific community. *Scientometrics*, 88:311–319.
- Mena-Chalco, J. P. and Cesar Junior, R. M. (2009). ScriptLattes: an open-source knowledge extraction system from the Lattes platform. *Journal of the Brazilian Computer Society*, 15:31 – 39.
- Mugnaini, R., Leite, P., and Leta, J. (2011). Fontes de informação para análise de internacionalização da produção científica brasileira. *PontodeAcesso*, 5(3).
- Newman, M. E. J. (2001). The structure of scientific collaboration networks. *Proc. of the National Academy of Sciences of the United States of America*, 98(2):pp. 404–409.
- Silva, F. M. and Smit, J. W. (2009). Organização da informação em sistemas eletrônicos abertos de Informação Científica & Tecnológica: análise da Plataforma Lattes. *Perspectivas em Ciência da Informação*, 14:77 – 98.