

DADOS GERAIS			
Título do projeto de pesquisa	Abordagem para integração automática de dados estruturados e não estruturados com uso de técnicas de aprendizagem em um contexto Big Data.		
Orientando	Keylla Ramos Saes		
Orientador (es)	Luciano Vieira de Araújo		
Semestre no curso, na data do workshop	<input type="checkbox"/> 2º semestre	<input type="checkbox"/> 3º semestre	<input checked="" type="checkbox"/> 4º semestre
Qualificação	<input checked="" type="checkbox"/> Qualificação já realizada em: 13/04/2018 <input type="checkbox"/> Realização da qualificação planejada para: dd/mm/aaaa		
Defesa	Prazo máximo para depósito: 03/09/2019 Realização da defesa planejada para: 10/2018		
Linha e Área de pesquisa	Gestão e Desenvolvimento de Sistemas: <input checked="" type="checkbox"/> BD <input type="checkbox"/> Engenharia de Software <input type="checkbox"/> Gestão de TI <input type="checkbox"/> IHC	Inteligência de Sistemas: <input type="checkbox"/> IA <input type="checkbox"/> Processamento Gráfico <input type="checkbox"/> Reconhecimento de Padrões	
Área de aplicação	<input type="checkbox"/> Ambientes corporativos <input type="checkbox"/> Bioinformática <input type="checkbox"/> Biometria <input type="checkbox"/> Dispositivos móveis <input type="checkbox"/> Educação	<input type="checkbox"/> Educação a distância <input type="checkbox"/> Governo eletrônico <input type="checkbox"/> Internet <input type="checkbox"/> Jogos <input type="checkbox"/> Jogos sérios	<input type="checkbox"/> Língua Natural <input type="checkbox"/> Linguística <input type="checkbox"/> Processos de Negócio <input type="checkbox"/> Quimioinformática <input checked="" type="checkbox"/> [Big Data] <input type="checkbox"/> Redes Sociais <input type="checkbox"/> Robótica <input type="checkbox"/> Saúde
Publicações associadas ao projeto de mestrado	<ol style="list-style-type: none"> 1) Integration of structured and unstructured data in a Big Data context: a survey (artigo em elaboração com envio acordado para ACM) 2) Approach to automatic integration of structured and unstructured data with use of learning techniques in the context of Big Data (artigo em elaboração com envio acordado para VLDB) 		

DESCRIÇÃO DO PROJETO DE PESQUISA

Contextualização / motivação	Essa pesquisa é parte de uma pesquisa maior sobre Data DNA.
Problema de pesquisa	Integrar dados com modelos de representação homogêneos não é um processo trivial, e a complexidade é ampliada ao considerar na integração modelos de representação de dados heterogêneos. Devido a variedade de modelos de representação, integrar dados de maneira manual pode tornar-se inviável. Como proposta para este desafio de integração, esse trabalho apresenta uma abordagem automática para apoiar a integração de dados com grande variedade de modelos de representação e alto volume de dados, ou seja, em um contexto Big Data.
Objetivo geral da pesquisa	Este estudo propõe uma abordagem para ampliação dos cenários de integração de dados automática entre modelos de representação de dados estruturados e não estruturados, como forma de apoio a geração de conhecimento, através de agilidade na disponibilização de dados integrados com redução de processos manuais.
Trabalhos relacionados	<ul style="list-style-type: none"> • Proposta de um processo de integração de dados heterogêneos utilizando uma virtualização de dados que possui um tradutor entre linguagem SQL e NOSQL. A proposta só aceita dados armazenados em SQL e dados armazenados no MongoDB, o que é uma limitação da solução. (LAWRENCE,2014). • Uma abordagem semiautomática propõe um mediador com algoritmos inteligentes com uso de Machine Learning para realizar a integração de dados contidos em esquemas distintos. A baixa acurácia da solução é uma limitação que será endereçada em trabalhos futuros. (DOAN; DOMINGOS, 2001) • O SEMINT é a proposição de um modelo de integração para dados estruturados com a utilização de algoritmos inteligentes. O trabalho apresenta como contribuição um processo semiautomático para integração de dados, porém limitados aos dados estruturados. (LI; CLIFTON, 2000)
Proposta para solução	(i) proposta uma nova abordagem;
Dados	Dados públicos extraídos de bases federadas como CKAN, Data One, e outras bases como Kaggle, e bases do governo, como bolsa família.
Validação	A validação da abordagem para integração de dados heterogêneos utilizará dados públicos extraídos de bases federadas como CKAN, Data One, e outras bases como Kaggle, e bases do governo para realizar integração de dados de maneira automática.

