

# Anotação de Trajetórias via Fusão com Trilhas de Mídias Sociais

Ricardo Gil Belther Nabo, Renato Fileto  
Cleto May, Lucas André de Alencar

<sup>1</sup>Departamento de Informática e Estatística (INE)  
Universidade Federal de Santa Catarina (UFSC)  
Florianópolis, Santa Catarina, Brazil

**Resumo.** *O aumento da utilização de dispositivos móveis tem ocasionado um grande crescimento na geração de trajetórias de objetos móveis. Entretanto somente as trajetórias muitas vezes não são suficientes para permitir a análise semântica dos movimentos. Junto ao crescimento da utilização de dispositivos móveis também aumentou a utilização de mídias sociais remotamente, onde postagens dos usuários podem ser vistas como rastros esparsos e anotados de seu movimento. Este trabalho propõe um método para fusão de trajetórias com dados provenientes de mídias sociais. Os resultados são coleções de trajetórias anotadas com texto inserido em mídias sociais. O método é implementado e avaliado em experimentos com trajetórias reais de postagens de usuários do Twitter, efetuadas na mesma região geográfica onde ocorreram as trajetórias.*

**Abstract.** *The increased use of mobile devices has led to a large increase in the moving objects trajectories generation. However, only the trajectories are often not sufficient to allow the movements semantic analysis. The increasing use of mobile devices has also increased the remotely social media use, where users' posts can be viewed as sparse and annotated traces of their movement. This paper proposes a method for fusing trajectories with data from social media. The results are collections of trajectories annotated with texts inserted in social media. The method is implemented and evaluated in experiments with real trajectories and postings of Twitter users, conducted in the same geographic region where the trajectories occurred.*

## 1. Introdução

A utilização de dispositivos móveis que permitem a coleta de coordenadas espaciais (eg., GPS, *smartphones*, *tablets*) tem tido um crescimento considerável nos últimos anos. Este crescimento tem acarretado a geração de grandes volumes de dados de trajetórias brutas.

Muitos trabalhos relacionados a mineração de padrões espaço-temporais vem sendo desenvolvidos na literatura. Dentre eles, estruturação de trajetórias [Spaccapietra et al. 2008, Xiu-li and Wei-xiang 2009] e anotação de trajetórias [Alvares et al. 2007, Yan et al. 2012]. A anotação de trajetórias é importante, pois somente os dados espaço-temporais das trajetórias geralmente não são suficientes para o enriquecimento semântico. As trajetórias brutas coletadas por dispositivos móveis (e.g., *smartphones*) quase sempre carecem de dados textuais. Por outro lado, dados que às vezes são coletados pelos mesmos dispositivos móveis e inseridos em mídias sociais (e.g.,

*tweets*, *posts* no Facebook) possuem informações textuais (e.g., comentários, *hashtags*) que podem ajudar a descrever e analisar semanticamente as trajetórias.

Este trabalho propõe um método para realizar a fusão de trajetórias brutas com *posts* de usuários em mídias sociais, utilizando como critério da fusão as coordenadas espaciais e os instantes de coleta de pontos de trajetórias. O método é validado utilizando uma base de dados de trajetórias brutas e dados da mídia social Twitter coletados no mesmo local.

O restante deste trabalho está organizado da seguinte maneira. A seção 2 define alguns conceitos fundamentais para a compreensão do trabalho. A seção 3 apresenta o método proposto. A seção 4 descreve e discute experimentos para validar o método. A seção 5 discute e compara este trabalho com outros trabalhos relacionados. Finalmente, a seção 6 apresenta as conclusões e trabalhos futuros.

## 2. Fundamentação

Esta seção apresenta alguns fundamentos e definições utilizados na descrição formal do problema abordado e do método de solução proposto neste artigo.

### 2.1. Trajetórias

Trajetoórias brutas são sequências temporalmente ordenadas de coordenadas espaço-temporais. Cada coordenada pode ser definida como um ponto.

**Definição 1. (Ponto espaço-temporal)** Coordenada espaço-temporal representada pela quádrupla:  $\mathbf{P}(Pid, x, y, t)$ , onde:

- $Pid$  é o identificador do ponto;
- $(x, y)$  é um par de coordenadas geográficas; e
- $t$  é um instante de tempo.

Um dispositivo móvel que coleta amostras de localizações na forma de pontos espaço-temporais dentro de um determinado intervalo de tempo gera uma trajetória bruta.

**Definição 2. (Trajetória Bruta - TB).** Sequência temporalmente ordenada de pontos,  $(Pid, x, y, t) (p_1, p_2, \dots, p_n)$  visitados por um objeto móvel, onde cada elemento desta sequência é representado pela tripla:  $\mathbf{RawTraj}(MOid, Tid, P_j)$ , onde:

- $MOid$  é o identificador do objeto móvel;
- $Tid$  é o identificador da trajetória; e
- $P$  é a referência para um ponto espaço-temporal (Definição 1).

Visando melhorar o desempenho e os resultados produzidos pelo processamento de TBs, seus pontos são agrupados de acordo com alguma característica em comum entre eles. Os grupos de pontos resultantes são denominados episódios.

**Definição 3. (Episódio).** Subsequência maximal de pontos de uma trajetória que satisfaçam um determinado predicado  $(P_{inicial} \dots P_{final}) : \implies \{true, false\}$ . Um episódio é representado pela quádrupla:  $\mathbf{Episódio}(Tid, Eid, EType, P_{inicial} \dots P_{final} (1 \leq inicial \leq final \leq n))$ , onde:

- $Tid$  é o identificador da trajetória a quem o episódio pertence;
- $Eid$  é o identificador do episódio;

- $Etype$  é o tipo do episódio (e.g. "stop", "move"); e
- $P_{inicial} \dots P_{final}$  é a subsequência de pontos que constituem o episódio.

Como a lista de pontos pertencentes a um episódio é uma subsequência maximal, somente referências para o primeiro e o último ponto de cada episódio precisam ser armazenadas em sua estrutura. Os demais pontos podem ser recuperados diretamente da TB, desse modo não duplicando informações.

Os episódios, quando temporalmente ordenados, geram outro tipo de trajetória, uma trajetória estruturada. Cada elemento da trajetória estruturada é um episódio.

**Definição 4. (Trajetória Estruturada - TE).** Sequência temporalmente ordenada de episódios não aninhados. Cada elemento da sequência é representado pelo par:  $StrTraj(STid, Ei)$ , onde:

- $STid$  é o identificador da trajetória estruturada; e
- $Ei$  é um episódio.

## 2.2. Dados de Movimentos colhidos em Mídias Sociais

Uma pegada de mídia social é o registro de uma interação entre um usuário e uma mídia social (e.g., Twitter, Facebook, Foursquare). Quando o usuário posta algo a informação associada (e.g., posição espaço-temporal, foto) fica gravada na respectiva mídia (e.g., Twitter, Facebook) e acessível via APIs específicas de cada mídia. Uma sequência temporalmente ordenada de pegadas constitui uma trilha.

**Definição 5. (Pegada de Mídia Social).** Registro em um sistema de mídia social de interação efetuada por um usuário, representado pela quintupla:  $SMF(MOid, SMFid, SMid, P, c)$ , onde:

- $MOid$  é o identificador do objeto móvel;
- $SMFid$  é o identificador da pegada;
- $SMid$  é o identificador da mídia social da pegada (e.g., Twitter, Facebook);
- $P$  é um ponto espaço-temporal como descrito na Definição 1;
- $c$  são os conteúdos das pegadas (e.g., tags, imagens, textos);

**Definição 6. (Trilha de Mídia Social - TMS).** Sequência temporalmente ordenada de pegadas de mídia social, do mesmo usuário. Cada elemento desta sequência é representado pela dupla:  $SMT(SMTid, SMF)$ , onde:

- $SMTid$  é o identificados da trilha de mídia social; e
- $SMF$  é a referência a uma pegada de mídia social.

Embora análogas em termos de estruturas de dados, TBs e TMSs são diferentes. Trajetórias usualmente têm melhor precisão espaço-temporal que trilhas. A amostragem de pontos de TBs usualmente é realizada em intervalos fixos e curtos (e.g., 10 segundos, 10 metros). Por outro lado, as postagens em mídias sociais são assíncronas (o usuário decide quando postar) e usualmente esparsas, além das TMSs possuírem anotações de usuários.

### 2.3. Descrição do Problema

Considere um conjunto de TBs (TB\_DB) e um conjunto de TMSs (TMS\_DB), como o ilustrado na Figura 1. Se considerarmos que alguns objetos moveis que geraram TBs (Definição 2) podem ser os mesmos (ou ao menos ter movimentos e intenções afins) àqueles que geraram TMSs (Definição 6), é possível fundir algumas trajetórias com trilhas. Entretanto para realizar essa fusão, é necessária a utilização de uma medida de correlação entre uma TB e uma TMS, tais como correlação espaço-temporal.

O desafio deste trabalho é desenvolver um método para fundir trajetórias provenientes de dispositivos GPS com TMSs. Em outras palavras, o problema é determinar pares trajetória-trilha que são espacial e temporalmente correlacionados. Por exemplo, na Figura 6, as trajetórias horizontais e verticais (representadas por linhas contínuas), podem corresponder às respectivas trilhas horizontais e verticais (representadas por linhas pontilhadas). Tal correlação pode indicar: (i) trajetória e trilha de cada par geradas pelo mesmo dispositivo; (ii) geradas pelo mesmo objeto móvel (e.g. pessoa) portando dispositivos distintos; (iii) pares trajetória-trilha com movimentos análogos (e.g., uma pessoa usando um *smartphone* para acessar mídia social de dentro de um veículo equipado com GPS).

O método proposto neste trabalho utiliza proximidade espaço-temporal entre episódios e pegadas, das respectivas trajetórias e trilhas, para medir tal correlação e calcular os coeficientes de correspondência de pares trajetória-trilha. Considera-se que a segmentação adequada das trajetórias foi previamente realizada para suportar a investigação das correspondências. A determinação de quais objetos móveis (referentes a trajetórias) correspondem a quais usuários de redes sociais (que geram trilhas) também está fora do escopo deste trabalho, sendo deixada para trabalhos futuros.

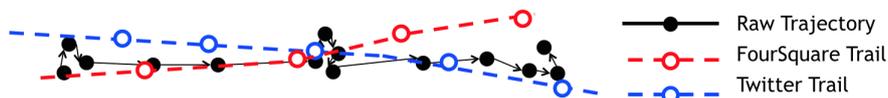


Figura 1. Ilustração do Problema

### 3. Método de Fusão Proposto

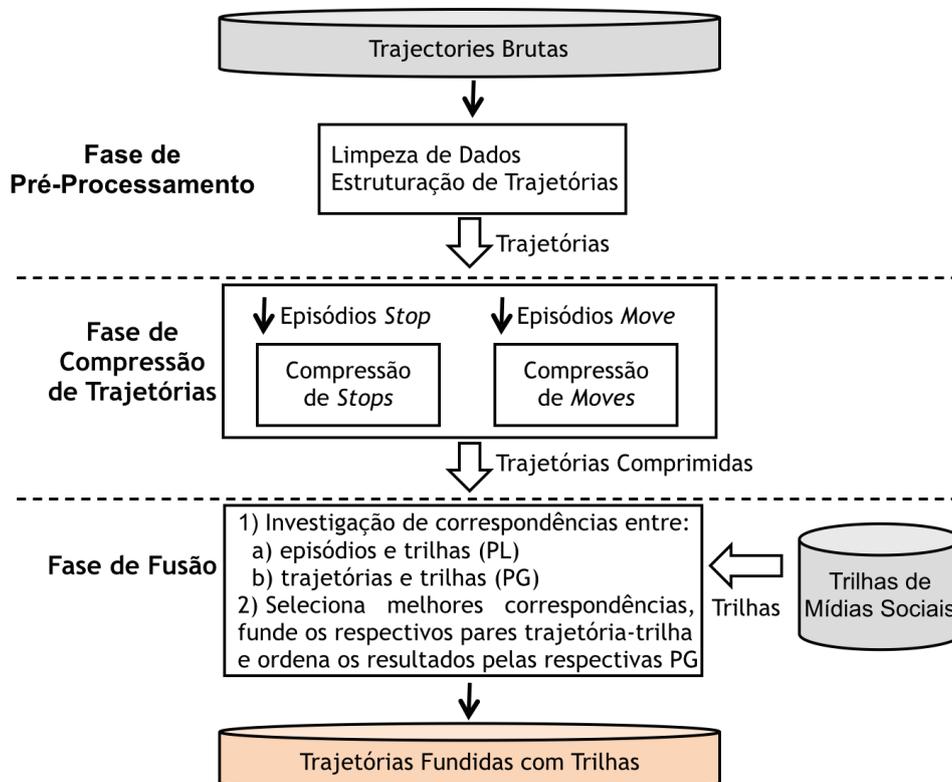
Esta seção apresenta um método para anotação de trajetórias mediante sua fusão com dados de mídias sociais. Inicialmente é apresentado o processo geral proposto, depois alguns detalhes de tarefas específicas deste processo.

#### 3.1. Um Processo para Fusão de Trajetórias com Trilhas de Mídias Sociais

O método proposto neste trabalho utiliza as definições descritas na seção 2. A figura 2 apresenta o processo proposto, o qual é composto de três fases: pré-processamento, compressão de trajetórias e fusão de trajetórias com trilhas. Todas as fases são flexíveis quanto aos métodos utilizados para sua implementação.

A fase de pré processamento das TBs consiste da limpeza e estruturação dessas trajetórias. A fase de compressão comprime os dados de trajetórias em uma representação que possa ser analisada de forma menos custosa que os dados brutos ou a mera agregação em episódios. A fase de fusão, foco principal deste trabalho, consiste na computação das probabilidades de correspondências entre pares trajetória-trilha. Os coeficientes de

correspondência local (CCL) de pares (episódio-trilha) espaço-temporalmente próximos são usados para computar o coeficiente de correspondência global (CCG) da respectiva trajetória com a trilha. Ao final, os pares (trajetória-trilha) com mais alto CCG são selecionados, fundidos e ordenados de acordo com seu respectivo CCG.



**Figura 2. Fases do método proposto**

### 3.1.1. Investigando Correspondências entre Episódios e Pegadas

O Coeficiente de Correspondência Local (CCL) é calculado para cada episódio de uma TE, verificando se há pegadas de mídias sociais dentro de um determinado *buffer* espaço-temporal ao redor de um episódio. Se há pegadas de uma trilha dentro deste *buffer* a chance da respectiva trilha e trajetória estarem correlacionada aumenta. A equação 1 apresenta o cálculo do CCL entre um episódio e as pegadas de uma mesma trilha que estão dentro do *buffer* espaço-temporal em torno do respectivo episódio.

$$CCL(Ep_i, SMT_j) = \frac{|SMF|}{|ALLFP|} \quad (1)$$

$Ep_i$  é o episódio,  $SMT_j$  é a trilha em que se está calculando o coeficiente, SMF é o conjunto de pegadas da  $SMT_j$  que estão dentro do *buffer* espaço-temporal em torno de  $Ep_i$  e ALLFP é o conjunto de pegadas de mídias sociais de toda a base de trajetórias que estão dentro do *buffer* espaço-temporal para este episódio.

A Figura 3 exemplifica o cálculo da CCL para duas TMSs ( $SMT_1$  e  $SMT_2$ ) e uma TE ( $StrTra_{j_1}$ ). As duas TMSs possuem pegadas dentro do *buffer* espaço-temporal



$CCL(Ep_1, SMT_1) = 0$  e  $CCL(Ep_1, SMT_1) = 0$ . O número de episódios de  $StrTraj_1$  (n) é 3. Portanto  $CCG(StrTraj_1, SMT_1) = \frac{\sum_{k=1}^3 P(Ep_k, SMT_1)}{3} = \frac{\frac{1}{2}+0+1}{3} = \frac{\frac{3}{2}}{3} = \frac{1}{2} = 0.5 = 50\%$ . A  $StrTraj_1$  tem 50% de chance de ser representada pela  $SMT_1$ .

Calculando do mesmo modo para a relação  $StrTraj_1 SMT_2$  temos:  $CCG(StrTraj_1, SMT_2) = \frac{\sum_{k=1}^3 P(Ep_k, SMT_2)}{3} = \frac{\frac{1}{2}+0+0}{3} = \frac{\frac{1}{2}}{3} = \frac{1}{6} \approx 0.17 \approx 17\%$ .

#### 4. Experimento

O método apresentado na seção anterior foi parcialmente implementado, com o intuito de validar a ideia de fusão de TBs e TMSs proposta neste trabalho. Para realizar a validação do método foram realizadas algumas simplificações referentes ao método, devido a sua alta complexidade de implementação. Dentro deste contexto o experimento descrito nesta seção utiliza somente *stops* para fundir trajetórias com TMSs.

Este experimento tem como objetivo comparar de forma qualitativa os tipos de anotações produzidas pelo método proposto neste trabalho e outros trabalhos da literatura, tais como [Yan et al. 2012] e [Alvares et al. 2007]. As entradas de dados descritas pelo método são duas bases de dados, uma de TBs e outra de TMSs coletadas na mesma região geográfica durante o mesmo período de tempo. Entretanto a construção de bases de dados de TBs e TMSs não é uma tarefa trivial, desse modo visando validar a fusão em si iremos utilizar bases de dados coletadas na mesma região, mas em um período de tempo diferente.

A base de dados de TBs foi selecionada de um grande conjunto de dados de táxis coletados na região de Fortaleza, durante o período de 20/07/2012 à 20/10/2012. Para este experimento foram selecionados quatro motoristas de táxis e foi considerada que cada trajetória teria a duração de 24 horas, gerando o total de 357 trajetórias, uma para cada dia de trabalho dos quatro taxistas. As TBs utilizadas neste experimento só consideram o caminho percorrido pelos taxistas, quando há um passageiro dentro do táxi. A base de dados de TBs é ilustrada na Figura 5.

A base de dados de TMSs consiste na aquisição de trilhas da mídia social Twitter provenientes de 227 usuários, na região de Fortaleza durante o período de (31/12/2013 à 02/01/2014). As trilhas possuem combinadas 1436 pegadas de mídia social, e não são segmentadas. Portanto o tamanho máximo de uma TMS neste experimento é de três dias. A base de dados de mídias sociais pode ser observada na Figura 6, os pontos azuis representam as pegadas da mídia social Twitter e as linhas lilases representam a ligação entre as pegadas geradas pelos mesmos usuários.

É responsabilidade da fase de pré-processamento garantir que as bases de dados estejam limpas e prontas para a aplicação do resto do método, também é nesta fase em que é realizada a estruturação de TBs. Neste experimento é utilizado o algoritmo o algoritmo *Cluster Based Stops and Moves on Trajectories* (CB-SMoT) [Xiu-li and Wei-xiang 2009], ele é responsável por encontrar episódios do tipo *stop* e *move*, considerando *clusters* de variação de velocidade e a densidade dos pontos.

Com a estruturação das TBs completas a fase de pré-processamento é encerrada, e é dado início a fase de compressão de trajetórias. Como já dito os episódios do tipo *Move* não foram considerados na fase de fusão e portanto não foram comprimidos durante a fase de compressão de trajetórias. Por outro lado os episódios do tipo *stop* precisam



**Figura 5. Base de dados de trajetórias brutas**



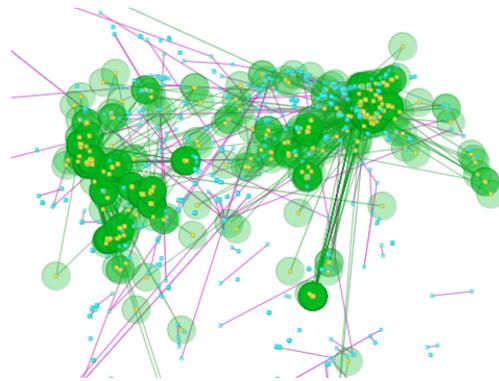
**Figura 6. Base de dados de trilhas de mídias sociais**

ser comprimidos, para tal neste trabalho iremos calcular o centroide de cada episódio do tipo *stop*. Desse modo cada *stop* se transforma em um ponto. Na Figura 7 está exposto a base de dados de TEs somente com os *stops* das trajetórias, já representados pelos seus centroides. Os pontos amarelos representam os centroides dos episódios do tipo *stop* e as linhas verdes representam a ligação entre os *stops* dos mesmos usuários.

Durante a fase de fusão é realizado o cálculo dos coeficientes de correlação locais e globais para todo o par trajetória-trilha das bases de dados, como descrito nas seções 3.1.1 e 3.1.2, entretanto como as bases de dados não são do mesmo período de tempo foi necessário realizar uma simplificação temporal para a fusão. Desse modo não foram utilizados parâmetros temporais para a realização dos cálculos das probabilidades locais e globais. A Figura 8 ilustra o cálculo das probabilidades locais e globais das bases de dados de TBs e de TMSs, a Figura também ilustra os *buffers* utilizados para o cálculo das probabilidades locais.



**Figura 7. Visualização de trajetórias estruturadas, utilizando os centroides de episódios *Stop***



**Figura 8. Trajetórias estruturadas com o *buffer* construído e trilhas de mídias sociais**

Após o cálculo dos coeficientes de correlação foi selecionada a TMS com a maior probabilidade global para todos os pares trajetória-trilha, e foi realizada a fusão destes pares, adicionando os dados inseridos pelo usuário no Twitter aos episódios do tipo *stop*. Esta fusão pode ser encarada como uma anotação de trajetórias.

#### **4.1. Resultados**

A Figura 14 apresenta um exemplo de TE e anotada produzida pelo método de fusão proposto, a partir dos dados usados nos experimentos (trajetórias de táxis e *tweets* em

Fortaleza). A trajetória, representada pela linha contínua foi segmentada em 3 episódios:  $Ep_1$ ,  $Ep_2$  e  $Ep_3$ . O método proposto detectou a correlação desta trajetória com a trilha composta pelas pegadas  $SMF_1$ ,  $SMF_2$ ,  $SMF_3$ ,  $SMF_4$ . Isso permitiu a associação das anotações a tais pegadas (listadas na tabela à direita) aos episódios espaço-temporalmente próximos das respectivas pegadas.

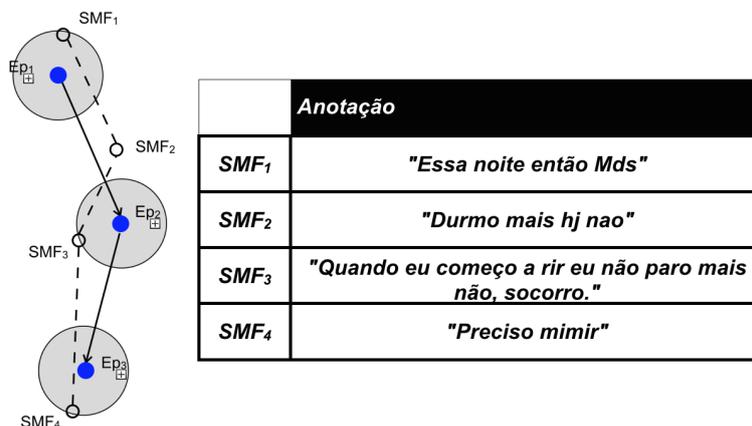


Figura 9. Exemplo de fusão

Os experimentos iniciais, relatados neste trabalho, sugerem a viabilidade de fundir trajetórias com TMSs, para obter anotações para as primeiras. Entretanto, é necessário a realização de experimentos com bases de dados do mesmo período de tempo para realizar tais afirmações.

## 5. Trabalhos Relacionados

O enriquecimento semântico de dados de trajetórias tem sido amplamente investigado na literatura por diversos projetos internacionais como por exemplo o MODAP (Mobility, Data Mining, and Privacy)<sup>1</sup> e SEEK (SEmantic Enrichment of trajectory Knowledge discovery)<sup>2</sup>. Diversos trabalhos dentro destes projetos buscam aumentar a quantidade de informações que se pode extrair dos mesmos mediante enriquecimento semântico [Alvares et al. 2007, Yan et al. 2012].

[Alvares et al. 2007] busca enriquecer as trajetórias utilizando o processamento de trajetórias (e.g., inferência do veículo que esta sendo utilizado pelo portador do objeto móvel baseado em sua velocidade), além de utilizar bases de dados de informações externas as trajetórias (e.g., pontos turísticos da região).

[Yan et al. 2012] propõe uma plataforma de anotação de trajetórias, esta plataforma consiste na estruturação de trajetórias em diferentes camadas, de algoritmos de processamento de trajetórias como citado em [Alvares et al. 2007] e utilização de bases de dados de pontos de interesse provenientes de *linked data*.

As anotações geradas por [Alvares et al. 2007] e [Yan et al. 2012] são fruto do processamento de padrões nas trajetórias e interpretação desses padrões, utilizando conhecimento externo as trajetórias, portanto as anotações não requerem somente uma base de dados para sua execução, mas também conhecimento especialista para sua anotação.

<sup>1</sup><http://www.modap.org/>

<sup>2</sup><http://www.seek-project.eu/>

O método aqui proposto, por outro lado, permite anotar trajetórias automaticamente, não sendo necessário conhecimento especialista. Além disso, até onde vai nosso conhecimento, esta abordagem baseada em fusão de trajetórias com TMSs é inédita na literatura sobre trajetórias de objetos móveis.

## 6. Conclusão e Trabalhos Futuros

Este trabalho propôs um método para fundir trajetórias com TMSs para agregar informações textuais a TEs. Esta fusão é realizada em três fases: pré-processamento, compressão de trajetórias e fusão de trajetórias com trilhas. A principal contribuição deste trabalho é a definição de um novo método para realização de fusão de trajetórias, entretanto para medir seu desempenho quantitativo comparado com outros métodos é necessário a realização de mais experimentos.

Tal método pode ser expandido para realização de fusões espaço-temporais entre bases de dados de TBS e trilhas de redes sociais, desde que as bases de dados pertençam ao mesmo período de tempo. As trajetórias anotadas geradas podem ser utilizadas para trabalhos futuros, como por exemplo o enriquecimento semântico de trajetórias com dados ligados [Fileto et al. 2013].

Como trabalhos futuros se espera a (i) realização de experimentos com base de dados pertencentes ao período de tempo; (ii) melhora da complexidade do algoritmo de fusão; (iii) enriquecer semanticamente trajetórias anotadas.

## 7. Agradecimentos

Este trabalho foi apoiado pelo projeto da União Europeia IRSES-SEEK, CNPq e CAPES.

## Referências

- Alvares, L. O., Bogorny, V., Kuijpers, B., de Macedo, J. A. F., Moelans, B., and Vaisman, A. (2007). A model for enriching trajectories with semantic geographical information. In *Proc. of the 15th Annual ACM International Symposium on Advances in Geographic Information Systems, GIS '07*, pages 22:1–22:8, New York, NY, USA. ACM.
- Fileto, R., Krüger, M., Pelekis, N., Theodoridis, Y., and Renso, C. (2013). Baquara: A holistic ontological framework for movement analysis using linked data. In Ng, W., Storey, V., and Trujillo, J., editors, *Conceptual Modeling*, volume 8217 of *Lecture Notes in Computer Science*, pages 342–355. Springer Berlin Heidelberg.
- Spaccapietra, S., Parent, C., Damiani, M. L., de Macedo, J. A., Porto, F., and Vangenot, C. (2008). A conceptual view on trajectories. volume 65, pages 126–146, Amsterdam, The Netherlands, The Netherlands. Elsevier Science Publishers B. V.
- Xiu-li, Z. and Wei-xiang, X. (2009). A clustering-based approach for discovering interesting places in a single trajectory. In *Intelligent Computation Technology and Automation, 2009. ICICTA '09. 2<sup>nd</sup> Int. Conf. on*, volume 3, pages 429–432.
- Yan, Z., Chakraborty, D., Parent, C., Spaccapietra, S., and Aberer, K. (2012). Semantic Trajectories: Mobility Data Computation and Annotation. volume 9, pages 39:1–39:34, New York. ACM Transactions on Intelligent Systems and Technology.