# Social-Sensor Composition for Scene Analysis

Tooba Aamir[1]([✉]), Hai Dong[1], and Athman Bouguettaya[2]

[1] School of Science, RMIT University, Melbourne, Australia
{tooba.aamir,hai.dong}@rmit.edu.au
[2] School of Information Technologies, The University of Sydney, Sydney, Australia
athman.bouguettaya@sydney.edu.au

**Abstract.** We consider the scene analysis as a service composition problem. A social-sensor cloud services composition model is proposed for the scene analysis. Our proposed model selects and composes social-sensor cloud services based on the user queries. Textual features of the social-sensor cloud services, i.e., description, comments, and meta-data of the social media images are used to reconstruct a scene. Our key contribution is an efficient and real-time composition of related images for scene analysis relying on meta-data and related posted information. Analytical results demonstrate the performance of the proposed model.

## 1 Introduction

The rich and explosive growth of social media data has resulted in the integration of social data into a range of data-centric applications [1,2]. Recent communication devices like smartphones, i.e., social-sensors provide the ability to embed sensor data directly into cloud-based social networks, i.e., social-sensor clouds [2,4]. Monitoring these social-sensors' activities provide multiple benefits in various domains. For example, urban management requires scene reconstruction and analysis in an area. Suppose the surveillance of the road segment through traditional sensors is limited in coverage. In such cases, social-sensors facilitate to fill in the information gap within events or happenings [5].

Social-sensor data, e.g., social media image meta-data and related posted data (e.g., location, description, and comments) are inherently multi-modal because of the different data formats and sources in those social media platforms. The multifaceted data poses a significant challenge for the efficient and real-time delivery of the social-sensors' data to the users [7,13]. In our previous work, we propose the *social-sensor cloud services* to provide an open, flexible, and reconfigurable platform for monitoring and controlling applications [4,5]. We abstract social-sensor cloud data, e.g., images' annotations (meta-data and related information like description and comments) as a service, i.e., *social-sensor cloud service*, to fulfill the users' information requirement [4,5].

This paper focuses on using the service paradigm as a vehicle to devise a method for scene reconstruction and analysis without carrying out actual

image processing. The aim is to provide the similar useful information about the required scene as image processing does [5]. A complete scene analysis needs images from multiple angles and different time intervals. In such cases, a composition of services is required to form multiple viewing angles to fulfill the users' requirement(s) [10,12]. In this regard, we have identified the following challenges:

– *Relevance Model for Spatio-temporal Cloud Service:* The accurate information regarding the context of the service is vital for better utilization and selection of the social-sensor cloud service as per user requirements [8]. The relevance of the service to the given query helps to ascertain whether the service is in the same context as of the query.
– *Spatio-temporal Composition.* Social-sensor cloud services composition becomes even more challenging in dynamic service environments characterized by changing conditions and context. An optimal composite service is a set of social-sensor cloud services, providing the best-suited services at any given time as per the users' query. A spatio-temporal composition aims to execute an optimal composition based on the functional attributes.

This paper accommodates the solution of the challenges mentioned above. We propose a composite service that will provide the user-required view and related information about any event or a happening for the scene analysis. The proposed composition model forms a tapestry in the spatial aspect and a storyboard in the temporal aspect. In the spatial aspect, the composition forms a scene by selecting images from un-coordinated users and placing them in a tapestry-like structure. In the temporal aspect, a timeline is formed by combining various tapestries to form the story of the event.

## 2   Motivation Scenario

Let us assume an accident occurred on 5th July 2016 around 8:30 pm, on the Pascovale Road, Glenroy. The crash involves two vehicles cars A and B crossing an intersection. The service user, i.e., the police has queried a scene analysis of the accident. The aim is to find the original behavior leading to the crash and the objects of interest, i.e., the vehicles or people involved. In such case, anyone in the area can act as a social-sensor by sharing images over a social network. We rely on these social-media images as social-sensor cloud services in the vicinity during that specific period to reconstruct the desired scene.

This work proposes a model for selection and composition of the social-sensor cloud services based on the user query. The query includes a region of interest, textual description and time of the queried event. The query includes (1) Query phrase, e.g., a car accident involving Car A and B on city-bound Pascovale Road. Car A and Car B are the objects of interest. (2) Query region, i.e., decimal longitude-latitude position. For example, *(−37.694264, 144.9131593)* covering the area of 10 m on all sides of the road. (3) Query time, e.g., 5th July 2016, from 8:25 pm to 8:40 pm.

The basic functional attributes of a social-sensor cloud service *Serv*, are abstracted from the social media image information. These include:

- Time $T$ of the service at which the image is taken.
- Description $D$ is a set of keywords or key-phrases providing additional information regarding the image, e.g., Car A crashes in Car B, Car accident.
- Location $L(x, y)$ is longitude and latitude position where the image is taken.
- Coverage $Cov$ of the image is defined as $VisD$, i.e., the maximum visible distance, covered by the image, $\overrightarrow{dir}$, i.e., the orientation angle of the image and $\alpha$, the angular extent of the scene covered by the image.

It is assumed that the available services are tagged with location and time. We index all the available services considering their spatio-temporal features using a 3D R-tree [4,6]. The search space is reduced by selecting the services that are spatio-temporally close to the querying location and contextually related to the query description. For example, at time $t_{-1}$, the descriptions of three images $img_1$, $img_2$ and $img_3$ show that Cars A and B were running along Pascovale Road city-bound, and Car C was taking the exit from M80 Ring Road. Cars A and B are objects of interest and therefore $img_1$, $img_2$ and $img_3$ are selected due to their contextual relevance to the query. Further, at time $t_0$, the description of an image $img_4$ shows that Car A stopped and avoided the collision with Car C. Therefore, Car C is considered as interacting with the object of interest Car A and $img_4$ is selected. Three images $img_5$, $img_6$ and $img_7$ in the spatio-temporal query region show that at the intersection, Car C ran the red light. At time $t_1$, four images $img_8$, $img_9$, $img_{10}$ and $img_{11}$ are selected due to their spatio-temporal and contextual relevance. Images $img_8$, $img_9$ and $img_{10}$'s description says that Car B and Car A crashed. Image $img_{11}$'s description says Car C escaped the accident scene.

11 services (images) are selected in this scenario. We cluster the selected services according to their spatio-temporal and contextual relationships. The contextual clustering is based on the interaction and relations between the services. The interactions and relationship between the objects of interest of the services are determined on the basis of the semantic similarity between the service description and the query description. The event-specific relationship describing the vocabulary dictionary provided by domain experts is used for this purpose. We assess the services for composability. The composability is assessed by predefined relations, explained in the relevance and composability models (Sect. 4). Finally, we build-up the composition, i.e., a visual summary by forming a tapestry-like scene. The composition is formulated by selecting the composable services covering the accident, the object of interests and the interacting object. The composition depicts the cars crashed and the cars involved, i.e., Car A and B crashed, and Car C escaped the crash scene.

## 3    Model for Social-Sensor Cloud Service

In this section, we have defined the social-sensor cloud service, selection, and composition model.

### 3.1   Model for an Atomic Social-Sensor Cloud Service

An atomic social-sensor cloud service *Serv* is defined by:

– *Serv_id* is a unique service id of the service provider *SocSen*.
– *F* is a set of functional properties of the service *Serv*.

### 3.2   Functional Model of an Atomic Social-Sensor Cloud Service

The functional requirements capture the intended behavior of an atomic service and form the baseline functionality. The minimal functional requirements associated with an atomic service and their information sources are:

– Social-sensor device: The basic functional attributes of a social-sensor cloud service associated with social-sensor device are time $t$, location $L(x,y)$ and coverage *Cov* of the sensor. We have discussed all these parameters in [5].
– Social-sensor service owner: Context *Con* of a social-sensor cloud service is associates with the service owner. It is the description of a service provided by the service owner. Context *Con* is defined by $D$ and $T$. Description $D$ of the service provides additional information regarding the image. It is assumed that the service's description includes complete detail of the service specifics related to the scene captured, e.g., objects captured, and their relations. Tags $T$ provide location and focus of the image.
– Social-sensor cloud: Interaction $I$ is the social network provided information regarding objects of interest in the image. It is assumed that the description includes detail of the objects of interests. This description is provided by the users of the cloud, i.e., social media, through comments. We assume that the information collected though comments is trustworthy. The trustworthiness of the comments is dealt in our previous work [3,9].

## 4   Social-Sensor Cloud Service Composability

In this section, we propose the social-sensor cloud service relevance and composability models for the social-sensor cloud services.

### 4.1   Model for Social-Sensor Cloud Service Relevance

The relevance of a service to a given query or another service helps to ascertain whether the service is in the same context as of the query or the other service. The relevance between two or more social-sensor cloud services can be described as spatio-temporal closeness, contextual relatedness and interaction relevance. The relevance between two services $Serv_1$ and $Serv_2$ can be defined as:

**Spatial Relevance.** $Rel_S$ means $Serv_1$ and $Serv_2$ are close in space boundaries and have similar coverage direction. This encompasses ($Serv_1.Cov_{(\alpha,dir)} \cong Serv_2.Cov_{(\alpha,dir)}$), i.e., similar in directions and angles AND $Serv_1.L =$

$Serv_2.L \pm \Delta$, i.e., close in the geo-location. Where, $\Delta$ is the max. allowed spatial difference.

**Temporal Relevance.** $Rel_t$ means $Serv_1$ and $Serv_2$ coincide in time, i.e., $((Serv_1.t_e = Serv_2.t_s \pm \varepsilon) \mid (Serv_1.t_s = Serv_2.t_s \pm \varepsilon) \mid (Serv_1.t_e = Serv_2.t_e) \pm \varepsilon)$. Where, $(Serv_1.t_e = Serv_2.t_s \pm \varepsilon)$ means the end time of $serv_1$ is close to the start time of $Serv_2$. $(Serv_1.t_s = Serv_2.t_s \pm \varepsilon)$ means the start time of $serv_1$ is close to the start time of $Serv_2$. $(Serv_1.t_e = Serv_2.t_e \pm \varepsilon)$ means the end time of $serv_1$ is close to the end time of $Serv_2$. $\varepsilon$ is the max. allowed time difference.

**Spatio-Temporal Relevance.** $Rel_{St}$ means $Serv_1$ and $Serv_2$ have overlap in time and space. This encompasses $Rel_S \cap Rel_t$.

**Contextual Relevance.** $Rel_C$ means $Serv_1$ and $Serv_2$ share same or almost similar context. This encompasses $(Serv_1.Con \cong Serv_2.Con)$. The contextual relevance is based on the textual similarity of the contextual descriptions of both services. Contextual relevance is calculated as a semantic distance between the descriptions of the services and the query [5]. Event specific relationships are used for the implementation of the similarity measure. These event specific relationships are described in the vocabulary dictionary provided by the domain experts. We have used $\theta$ to define $related_{LIN}(Serv_1.Con, Serv_2.Con)$. The higher value of $\theta$ shows higher similarity in context.

**Interaction Relevance.** $Rel_I$ means $Serv_1$ and $Serv_2$ both share objects of interest in the coverage (refer Sect. 4.1). This encompasses $(Serv_1.I \cap Serv_2.I)$

### 4.2   Model for Social-Sensor Cloud Service Composability

The spatio-temporal and contextual composability of two or more social-sensor cloud services can be defined as four instances:

- $(Rel_{St} \cap Rel_C)$. Two or more services are composable if these services are spatio-temporally and contextually relevant.
- $(Rel_t \cap Rel_C)$. Two or more services are composable if these services are temporally and contextually relevant. In such cases, services might be located outside the region of interest but still capture a scene inside.
- $(Rel_S \cap Rel_C)$. Two or more services are composable if these services are spatially comparable and contextually relevant. In such cases, services are available either before or after the required period.
- $(Rel_C \cap Rel_I)$.Two or more services might be composable if these services share context and objects of interest. In such cases, services might be located outside the region of interest but still capture some related objects of interest.

## 5   Social-Sensor Cloud Service Composition Approach

We propose an approach to filter, select and compose the best available social-sensor cloud service to form a visual summary according to the user's query. The

composition is achieved by constructing the information context of the service with the functional. The composite service comprises a set of selected atomic services to form a visual summary of the queried event. The visual summary offers an arrangement of the 2D images, forming a tapestry-like scene of the required event. Our approach aims to efficiently compose the available services into a single composite service that matches with the users' requirements.

A query $q$ can be defined as $q = (Rgn, des, t_s, t_e)$, giving the region of interest, description and time of the required service(s).

– $Rgn = \{P < x, y >, l, w\}$ [5], where $P$ is a geospatial co-ordinate set, i.e., decimal longitude-latitude position and $l$ and $w$ are length and width distance from $P$ to the edge of region of interest.
– $t_s$ is the start time and $t_e$ is the end time of the query.
– $des$ is a phrase describing the query. Query description includes details of the objects of interests $obj$, i.e., objects involved and the context of the query $cont$, i.e., the scene to be captured.

### 5.1   Social-Sensor Cloud Service Selection

The indexing and spatio-temporal filtering of the services enable the fast discovery of the services. We index all the available services using a 3D R-tree [4] and select the services inside the bounded region of interest [5]. Next, the services are selected and classified based on the relevance between the services, the queried scene and the objects of interest. It might happen that the service does lie spatio-temporally in the query area $Rgn$, but has no contextual relation with the query $q$ or has too much noise concerning unwanted information. In such cases, the object(s) of interest and behavior relations are used for the service filtration. The contextual relevance of all the services to a query's scene and objects of interest are assessed. Using previous research as reference we have set the value of threshold $\theta = 0.5$ for the contextual relevance [14]. The services related to the queried scene and objects of interests are selected. The services are classified in three sets according to their relevance: (1) spatio-temporally and contextually relevant services $S_{StC}$, (2) spatio-temporally relevant and interacting services $S_{StI}$ and (3) contextually relevant and interacting services $S_{CI}$.

### 5.2   Social-Sensor Cloud Service Composability Assessment

The composability rules aims to construct a composite service. Composability assessment among component services is based on their spatio-temporal and contextual parameters. The relevance and overlap is considered to define the composability relations between the services, e.g., $Serv_1$ and $Serv_2$. We aim to define composability of the service as quantitative relations. The relevance between the services is an arithmetic mean of the considered parameters. It is calculated as:

$$Rel(Serv_1, Serv_2) = [(Rel_{St}(Serv_1, Serv_2)+ \\ Rel_C(Serv_1, Serv_2) + Rel_I(Serv_1, Serv_2))] \tag{1}$$

where, $Rel_{St}(Serv_1, Serv_2)$ is based on the time of the services and their proximity in space. $\lambda$ is the shortest distance between $Serv_1$ and $Serv_2$ and $\vartheta$ is the difference between coverage angles $Serv_1.Cov_{dir}$ and $Serv_2.Cov_{dir}$. The thresholds for the spatial relevance are set as $\lambda_{thr}$ for distance and $\vartheta_{thr}$ for $\overrightarrow{dir}$. Therefore, the services are considered spatio-temporally relevant if difference between the distance and direction of the services is below the threshold. $Rel_C(Serv_1, Serv_2)$ is the semantic distance between the descriptions of $Serv_1$ and $Serv_2$ (Refer Sect. 5.1). $Rel_I(Serv_1, Serv_2)$ is the count of the mutual objects of interest in $Serv_1$ and $Serv_2$. The overlap between the services is considered:

$$Overlap(Serv_1, Serv_2) = Overlap_{spatial}(Serv_1, Serv_2) \qquad (2)$$

The quantitative value of the mutual composability is calculated as:

$$Comp(Serv_1, Serv_2) = Rel(Serv_1, Serv_2) - Overlap(Serv_1, Serv_2) \qquad (3)$$

A geographic coverage patch $GeoPatch$ is formed to assess the composability of each service from the spatio-temporal and contextual selection $S_{StC}$. A set $N$ of the spatio-temporally nearest services is selected for each $GeoPatch$. The mutual composability $Comp$ is calculated with each service in $N$. The process of calculating the mutual composability of the services is repeated with the sets $N'$ and $N''$. $N'$ is the set of the nearest services concerning the spatio-contextual and temporal-contextual relevance. $N''$ is a set of the nearest services based on the contextual relevance and interaction. The assessment process of the mutual composability is based on relevance and overlap of the services.

### 5.3    Social-Sensor Cloud Service Composition

The composition is handled as sewing a tapestry to form the scene. We start with the central piece, concerning space and time, and build a tapestry around it. The build-up is based on selecting the best composable services from the set of nearest services. The best neighbor service is with the maximum relevance and the minimum overlap.

The composition covers the visual summary of the whole queried *scene*, i.e., all objects of interest and their context. We choose the central service $Serv_c$ in terms of space and time from the spatio-temporal and contextual selection. We further add $Serv_c$'s neighbors to a separate pool. We assume that the central service is in the middle of the spatio-temporal dimension. Next, we extract the best neighbor service $Serv_{k.bn}$ from the pool and place it with $Serv_c$ by joining the patch. $Serv_{k.bn}$ is selected according to the maximum composability. We add neighbors of $Serv_{k.bn}$ to the pool. The process of selecting the best neighbor and joining to patch continues until we have any service in the pool. We reassess the composability of the remaining services and start again with the nearest service if the pool is empty. Spatial gaps in the composition are assessed after the utilization of all services from the spatio-temporal and contextual selection. $Comp.C$ is the total coverage of the services in the composition overlapping the

bounded region *Rgn* and within time t$_s$ and t$_e$. The relationship between *Serv* and *q.Rgn* can be illustrated as:

$$Composition \longrightarrow \{Comp \in \cup_{i=1}^{n} Serv \mid (Comp.C \cap Q.Rgn)\cap$$
$$Rel_C \cap Rel_I, t_s \leq t \leq t\_e\} \tag{4}$$

In our previous work, we have discussed the coverage of the composition and gap assessment [5]. We estimate and select an arbitrary neighbor $Serv_{kc'}$ if there are any spatial gaps. Next, the best nearest service $Serv_{k.bn}$ from the set of the spatio-temporally relevant and interacting services. The process of selecting and joining the services continues until we fill in the gaps and get the maximum available coverage. The composite service is a series of spatial tapestries in time, providing a timeline of the visual summary of the event.

## 6    Experiment and Evaluation

We focus on evaluating the proposed approach using the real dataset. The set is a collection of 10000 user uploaded images downloaded from social networks (flicker, twitter, google+). We had extracted their geo-tagged locations, the time when an image was captured, post description and tags to create the services. Further, the camera direction $\overrightarrow{dir}$, the maximum visible distance of the image *VisD* and the viewable angle $\alpha$ are abstracted as the functional property *Cov*.

We generated eight different queries based on the locations and events in our dataset. We have evaluated the service composition based on the spatial relevance in the first part of the experiment. The result of these experiments is evaluated upon the traditional image processing technique SIFT (Scale-Invariant Feature Transform) [11]. We used images' geolocation information, associated directions and viewing angles to gather an associated image dataset *I* from Google Street View of the area of interest *R*. We first downloaded 360° views of Google Street View using GPS from the image and collected the views related to the service. Further, we compared the similarity between images in the composition and the image set *I* by SIFT features. This comparison is achieved by individually comparing the key point feature vector of the images in *I* and images in the composition, and finding the images' matching features based on the Euclidean distance of their feature vectors. Further, we assessed if the images in the composition are correctly positioned in spatial relations. The evaluation of the similarity threshold is set around 60%. 40% noise margin is given due to traffic and pedestrian obstruction in the images.

We have assessed how useful the composite service is in completing the contextual storyboard in the second part of the experiment. The assessment is done by manually analyzing the composition for the spatial-temporal and contextual coverage. The effectiveness of the composite service is assessed upon the selection and composition of the related and accurate services. It is assessed if the composite service contains the required object(s) of interest and their behavior according to the user query.

**Table 1.** Relative accuracy in spatio-temporal coverage

|  | Queries themes | | Accuracy rate of SIFT selection |
|---|---|---|---|
| Event-oriented queries | Q1 | Bourke street accident | 57.9% |
| | Q2 | F1 race, Melbourne | 54.5% |
| | Q3 | Essendon airport crash | 36.8% |
| | Q4 | CBD random accident | 54.5% |
| Location-oriented queries | Q5 | Melbourne night | 81.5% |
| | Q6 | Melbourne central | 78.6% |
| | Q7 | Melbourne trams | 75.0% |
| | Q8 | Elishbeth street, Melbourne CBD | 69.2% |
| Average | | | 63.5% |

**Table 2.** Precision and recall

|  | Queries themes | | Precision | Recall |
|---|---|---|---|---|
| Event-oriented queries | Q1 | Bourke street accident | 51% | 72% |
| | Q2 | F1 race, Melbourne | 55% | 88% |
| | Q3 | Essendon airport crash | 73% | 79% |
| | Q4 | CBD random accident | 77% | 66% |
| | **Average** | | **64%** | **76%** |
| Location-oriented queries | **Q5** | Melbourne night | 88% | 79% |
| | Q6 | Melbourne central | 68% | 55% |
| | Q7 | Melbourne trams | 80% | 53% |
| | Q8 | Elishbeth street, Melbourne CBD | 80% | 74% |
| | **Average** | | **79%** | **65%** |

## 6.1 Evaluation

We have evaluated the proposed approach by (1) accuracy in the spatial coverage of the user required region, (2) effectiveness in selecting the related services (precision), and (3) effectiveness of the composite service in capturing the required context, i.e., the object(s) of interest and their behaviors (recall). All images and the composed services are manually analyzed by a human to form a baseline.

We have assessed the composite services by comparing the similarity between the service image and the Google street view. SIFT image processing is used for the comparison of all the eight queries (Table 1). We observed that approximately 63% of services in the compositions are accurately categorized in space. The 37% error rate was reasonable due to the noise in the images. Noise is an obstruction in the image affecting the scene building. For example, a vehicle obstructing the building of interest can be considered as noise. Further, we have assessed the composite services by manually analyzing the effectiveness of selecting the relevant spatio-temporal services, i.e., precision (Table 2). The average precision of the proposed approach for the location-based queries is 78% and for the event-based queries is 64%. The effective spatio-temporal and contextual coverage are assessed by recall (Table 2). The average recall of the proposed approach for the location-based queries is 65% and for the event-based queries is 76%. The results show that the values of precision are higher for the location-oriented queries, e.g., Melbourne Central (Q6). The values of recall are higher for the event or scene-oriented queries, e.g., Bourke Street Accident (Q1). Therefore, it

is concluded that our proposed approach effectively helps in the accurate composition of the services for the scene analysis. The proposed approach considers the related contextual data that describes the situation from various aspects, e.g., what has happened, where it happened, who is involved and what the effects on surrounding area.

## 7   Conclusion

We propose a social-sensor cloud service composition approach based on the spatio-temporal and contextual relevance. Our experiments evaluate the proposed approach for an accurate and effective composition. We plan to focus on the optimal social-sensor cloud service composition based on the uncertain time, location and context requirements.

## References

1. Rosi, A., Mamei, M., Zambonelli, F., et al.: Social sensors and pervasive services: approaches and perspectives. In: Proceedings of PERCOM (2011)
2. Aggarwal, C.C., Abdelzaher, T.: Social sensing. In: Aggarwal, C.C. (ed.) Managing and Mining Sensor data. Springer, Boston (2013). https://doi.org/10.1007/978-1-4614-6309-2_9
3. Aamir, T., Dong, H., Bouguettaya, A.: Trust in social-sensor cloud service. In: Proceedings of IEEE ICWS (2018)
4. Aamir, T., Bouguettaya, A., Dong, H., et al.: Social-sensor cloud service selection. In: Proceedings of IEEE ICWS (2017)
5. Aamir, T., Bouguettaya, A., Dong, H., Mistry, S., Erradi, A.: Social-sensor cloud service for scene reconstruction. In: Maximilien, M., Vallecillo, A., Wang, J., Oriol, M. (eds.) ICSOC 2017. LNCS, vol. 10601, pp. 37–52. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-69035-3_3
6. Neiat, A.G., Bouguettaya, A., Sellis, T., Ye, Z.: Spatio-temporal composition of sensor cloud services. In: ICWS (2014)
7. Bouguettaya, A., Singh, M., et al.: A service computing manifesto: the next 10 years. In: CACM (2017)
8. Wang, H., Shi, Y., et al.: Web service classification using support vector machine. In: Proceedings of IEEE ICTAI (2010)
9. Aamir, T., Dong, H., Bouguettaya, A.: Stance and credibility based trust in social-sensor cloud service. In: Proceedings of WISE (2018)
10. Ghari Neiat, A., Bouguettaya, A., Sellis, T.: Spatio-temporal composition of crowd-sourced services. In: Barros, A., Grigori, D., Narendra, N.C., Dam, H.K. (eds.) ICSOC 2015. LNCS, vol. 9435, pp. 373–382. Springer, Heidelberg (2015). https://doi.org/10.1007/978-3-662-48616-0_26
11. Lowe, D.G.: Distinctive image features from scale-invariant key points. IJCV **60**, 91–110 (2004)

12. Li, L., Liu, D., Bouguettaya, A.: Semantic based aspect-oriented programming for context-aware web service composition. Inf. Syst. **36**(3), 551–564 (2011)
13. Bouguettaya, A., Nepal, S., et al.: End-to-end service support for mashups. In: IEEE TSC (2010)
14. Mihalcea, R., et al.: Corpus-based and knowledge-based measures of text semantic similarity. In: Proceedings of AAAI (2006)