

# **An initiative towards a simplified international in-depth accident database**

D Ockel, J Bakker and R Schöneburg

Daimler AG, 71059 Sindelfingen, Germany

**Abstract** - While accident statistics on a national level are provided by many countries, there is a need for international data that includes more detailed information about the accident, so called in-depth data. As a consequence, accident data projects have been emerging in different regions of the world. This creates a need for comparable and mergeable data from different countries, enabling the use of already existing accident data resources and helping to expedite the improvement of global road safety.

While existing approaches focus that mostly on building a comprehensive accident database from scratch, the iGLAD project (Initiative for the Global Harmonization of Accident Data) attempts a more pragmatic approach by building on top of the work already accomplished in this area and complementing it. The target of iGLAD is to help setting up an additional dataset as a compatibility layer between already existing world wide data sets and integrating the structure of these by defining a common data scheme. This dataset is limited to the common denominator between the existing data sets and is inherently rather small and simple. Eventually, an individual converter for each participating accident investigation group will be built that enables pooling all data sets in a common repository. This not only saves costs and time, and hence makes such a target more feasible, but also creates data that is usable right from the start.

This paper gives an overview of the current status of iGLAD and first steps taken. Additionally, some methodological aspects are discussed, next to a glance at other projects working currently on related issues, providing additional input for iGLAD. Finally, an overview of next steps and intended future work is given.

## **1. INTRODUCTION**

While there is a strong demand for international accident data providing more detail than most national record based schemes, building up such a data pool is a great challenge and has not been achieved to date. Even if limited to EU-27, a centralized multinational in-depth accident data project is huge in terms of organizational effort and financial demand and thus hard to realize. The iGLAD project (Initiative for the Global Harmonization of Accident Data) follows a different approach which is bottom-up and of a more pragmatic and evolutionary character. Starting with different but already available pieces of data which, put together as they are, make up an inhomogeneous data set at first, iGLAD strives to build a usable and more homogeneous data set out of it. Furthermore, in the long run iGLAD strives for the convergence of in-depth accident data sources, as more and more data becomes globally available.

## **2. HISTORY**

Being a young project that kicked off at the end of 2011, iGLAD's history is rather brief. Basic discussions started within the GIDAS (German In-Depth Accident Study) [1] steering committee, triggered by requests from emerging in-depth data projects in other countries seeking support and best practices in how to set up in-depth data investigations. Central point of a detailed in-depth investigation is the code book which essentially reflects the complete data scheme. While there are differences between countries, e.g. in infrastructure and car fleet, the core structure of such an accident data scheme and the needs of the data users from different organizations (governmental, automobile industry OEMs and suppliers, educational and research institutes) are very similar.

Experiences with GIDAS and other accident investigation projects show that a full scale in-depth project can become very complex and challenging to maintain. There has been agreement within the GIDAS administration to share experiences and best practices of GIDAS's well-proven data scheme to facilitate comparability with emerging in-depth projects. In this spirit, iGLAD was initiated as a working group at the FIA Mobility Group in October 2010 and will address this challenge. Supported by FIA and ACEA, the goal of group is to define a common standardized accident data set as an effective foundation for developing and measuring road safety policy endorsements and interventions.

It shall also establish how this data set helps to achieve the goals of the “European Road Safety Action Programme” [2] and the „Decade of Action for Road Safety“ [3].

iGLAD was confirmed by the FIA Manufacturers Commission in March 2011. After presenting the basic concepts of iGLAD to NHTSA/NCSA, especially the NASS group in April 2011 and at the VDI congress [4], the project kick-off meeting followed on 30 September 2011 at ACEA, also marking the beginning of common and cooperative tasks of FIA and ACEA within the iGLAD project. One such task is a project assigned by FIA to analyse the traffic safety data situation in low-income and emerging countries, complementing the efforts of ACEA which initially address in-depth projects in higher and middle-income countries. The first iGLAD working group meeting in March 2012 comprised a more detailed discussion on the common data scheme and steps necessary for a standardized data set. The next section expands on the results.

### 3. STATUS

The kick-off meeting in September 2011 focused on congregating the representatives of the different projects to exchange the principal willingness for co-operation and establish which data could be available for use in a common database.

Participating organizations of iGLAD currently are: IRTAD (OECD), MUARC (Australia), VSRC (UK), DaCoTA (EU), Renault and LAB (France), ARU, VUFO GmbH and GIDAS (Germany), SAFER (Sweden), Graz University (Austria), MTI (Poland), NHTSA (USA), George Washington University (USA), Applus IDIADA Group (Spain), FIA and ACEA (Belgium). Additionally, the following organizations have been kept in the loop: NCSA, IIHS (USA), EC, ERSO (EU), SRA (Sweden), JP Research (India, USA), Uni Pavia (Italy), ITARDA (Japan), KATRI (Korea), CDV (Czech Republic), and the “Kuratorium für Verkehrssicherheit” (Austria).

	BAST	LAB	IDIADA SP	IDIADA CZ	ITS	JP Research	Uni Pavia	OTS	CCIS	DaCoTA
<b>General Data</b>										
1. Date, place	yes	yes	yes	yes	yes	yes	yes	yes (...)	yes (...)	yes (...)
2. Original police recorded data (for weighting)	yes	yes	police agreement	police agreement	yes	yes (...)	yes	yes	yes	yes
3. Accident description	yes	yes	yes	yes	no	yes	yes	yes	yes	yes
4. Pictures	yes	yes	yes	yes	no	yes	yes	yes	yes	yes
5. Accident sketch (scaled) with final positions and objects	yes	yes	police agreement	yes	no	yes	yes	yes	no	yes
6. Accident and collision type	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes
7. Environment: type of street, light & weather conditions, urban/rural	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes
8. Emergency arrival (time)	yes	no	police agreement	yes	no	yes (...)	no	yes	no	yes
<b>Vehicle / Pedestrian</b>										
1. Type of vehicle (make, model)	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes
2. Registration year	yes	yes	yes	yes	possible	yes	yes	yes	yes	yes
3. Vehicle parameters (mass, engine type, number of seats, shape)	yes	yes	yes	yes	possible	yes	yes	yes	yes	yes
4. Deformation (VDI/CDC)	yes	yes	yes	yes	possible	yes	yes	yes	yes	yes
5. Systems (passive/active safety)	yes	yes	integral safety	no	possible	yes	yes	yes	yes	yes
6. Pedestrian information (if any involved in accident)	yes	yes	yes	yes	yes	yes	yes	yes	No peds	yes
<b>Occupant</b>										
1. Age, gender, weight, height	yes	yes	certain	yes	age, gender	yes (...)	yes	yes	yes	yes
2. Injury severity (MAIS), AIS of body region, fatal injury	yes	yes	yes	yes	possible	yes (...)	yes	yes	yes	yes
3. Restraint systems (presence, use, and deployment)	yes	yes	yes	yes	no	yes	yes	yes	yes	yes
<b>Reconstruction</b>										
1. Collision opponent	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes
2. Collision speed	yes	yes (...)	yes	certain	no	yes (...)	yes	yes	no	yes
3. Driving speed	yes	yes (...)	yes	yes	no	yes (...)	yes	yes	no	yes
4. Delta-v	yes	yes	yes	certain	no	yes (...)	yes	yes	yes	yes
5. EES	yes	yes	yes	certain	no	yes (...)	yes	yes	yes	yes
6. Variables above for each collision (primary, most severe)	yes	yes (...)	yes	yes	no	yes	yes	yes	yes	yes

Table 1 – Overview of data availability survey for first draft of common data scheme (excerpt).

The participating organizations presented details about their relationship, support and possible contribution to iGLAD. A first impression of the potential of a common database could be derived from a survey prepared beforehand that consisted of a template reflecting a first and very rough draft of a common data scheme. The templates were completed by each participating organization

according to their own in-depth data studies. A selected overview of the first draft common data scheme and availability in the participating projects is shown in table 1, where ellipsis denote restricted availability for this particular item, and details are omitted. It can be concluded that a common data set comprising the data of the currently participating organizations already has sufficient potential and encourages further work in this regard.

More details of the common data scheme were discussed at the first working group meeting in March 2012. Experience from the DaCoTA project helped to identify variables in the data scheme that need special attention in the context of internationalization. Also, a study on converting data from GIDAS to DaCoTA showed that these two data sets can be mapped onto each other, at least at the core level which affects the domain of the iGLAD data set. It was agreed to start a pilot study where each data supplier converts a small set of accidents into the current version of the common data scheme. This should show the feasibility of the approach and give an indication about the resulting data set that could be provided by the iGLAD project. Confidentiality issues and modalities for the exchange of data need further discussion and a subgroup was created to handle these in more detail. A second subgroup will consider future co-operation with new data investigations in Europe, especially in the context of the DaCoTA project.

#### **4. METHODOLOGICAL ASPECTS**

This section addresses related projects and how iGLAD differs from them. Also, principal methodological issues in context of the iGLAD approach are shown with possible solutions to handle them.

##### **Related projects**

While other projects have already addressed the need for a multinational in-depth accident database, these mostly have a different focus than iGLAD. Even more, iGLAD should be able to contribute to and improve the current situation by complementing the work previously done and not replacing it. This should lead to an overall more complete solution which should be beneficial for all involved parties. The following is an (uncomplete) overview of what has been done or still is underway in this field.

Firstly, there are projects that define an in-depth data scheme or standard themselves and some of them also generate data or set up own teams to collect data: STAIRS [5], EACS [6], Pendant [7], TRACE [8], SafetyNet [9] and the currently running project DaCoTa [10].

Then, there are projects with a special focus: motorcycle study MAIDS [11] and truck study ETAC [12].

Also very important in this context are national statistics on a macroscopic level. They contain basic accident numbers of a larger scale of different countries and take care of a harmonized understanding and definition of the parameters contained (e.g. the definition of fatality): IRTAD [13] and CARE [14].

CARE data is based on disaggregated data and thus has access to individual accidents, but is limited to EU countries. IRTAD also contains non European countries.

Finally, CADaS [15] has introduced a reduced data scheme and a proposal to gradually implement it on a national level in Europe. This scheme contains 73 variables and is mainly based on CARE. The detail of data is somewhere in between national level and in-depth, reconstruction information is not planned for inclusion. Of the projects listed here, CADaS comes closest to what is intended within iGLAD.

However, the main differences of these projects compared to iGLAD are that most of them are limited to EU countries, they have a fixed time frame (except for the national data projects IRTAD and CARE) and some of them have a special focus. All of the projects either investigate data for themselves or design a data scheme to be filled by future projects, which can be considered as a top-down approach.

IGLAD follows an alternative bottom-up approach by employing what is already there. The basic difference is that no accident investigation teams are installed and no new accident case data is created within iGLAD; it is rather intended to provide the “glue” between existing projects. Also, iGLAD is designed to be as simple as possible, leaving complex details to the individual projects. This not only has practical reasons, since a very fine grained standard containing a long list of parameters can hardly be a basis for a common data set. Also, with a simple standard, details and country specific issues are left to the particular in-depth study, complementing these studies and bringing them closer to a wider and more globally oriented audience of researchers.

### **Define the common data scheme**

The approach taken by iGLAD is very pragmatic: See what is already there and build on top of it. Also, the result must be kept small and simple. iGLAD strives find an optimum between unifying a limited number of parameters and maintaining realistic targets and effectiveness. To achieve this, the different interests of the supporting members need to be carefully balanced. The level of detail provided by the resulting common data subset is not only a technical question, but also depends on the interests of the consortium partners. The result should be a well-balanced data set, where each party provides and receives comparable value. As an additional benefit for the data suppliers, the common data subset might spawn interest for further analyses (or contracted analyses) of their detailed data, i.e. the data available beyond that provided by the common subset.

Nevertheless, despite its target on simplicity, it is important that the data creates a useful basis for typical accident data analysis questions. To accomplish this, the working group needs to prepare relevant use cases of the data for demonstration purposes.

### **Adapt different samples**

One crucial step for bringing the iGLAD concept to life is adapting the different accident data samples to form a best possible homogeneous data set, resulting in a quality and expressiveness of the data that is sufficient for real life problems.

The issue of adapting the data samples has two largely independent dimensions. First, on a parameter level, data converters have to be built that are able to map the different schemes to the common data scheme. Second, the different sample characteristics must be compensated for differences or bias, which only affects the case level. The good thing is these two issues can be handled separately: Case and parameter level can be considered as orthogonal, i.e. they can vary independently from each other.

The first issue is mainly a work that has to be done only once by setting up a data converter for each sample. Close knowledge of the parameters in each sample are important for defining the most appropriate mapping between the values. Simple example: it is easy to convert between units like inch and cm, but mapping two different accident types involves accounting for regional and systematic differences.

The second issue can be addressed by the use of multinational statistics like IRTAD and CARE which can serve as a link between the in-depth data samples, provided that the in-depth samples include some parameters of the national statistics and that they are large enough to reflect the real world accident situation in the specific country.

IRTAD already does a lot of work to harmonize national statistics and supports iGLAD. Adapting a sample to national statistics can be accomplished by introducing weighting factors to compensate for differences in the sample characteristics of in-depth and national data.

Assuming that the raw data from national statistics can be accessed on a case level, arbitrary multidimensional tables of all available parameters can be generated easily for a specific country. In this ideal case, weighting factors can be calculated directly for the parameters that are present in both the national statistics and the in-depth sample. These weighting factors should be updated regularly (e.g. yearly) to account for changes over time. However, in most cases, accessibility of national data is more restricted and only distributions of single parameters or crosstabs of two or three parameters are available. Then, weighting factors can be calculated by filling up the contingency table of the weighting parameters with an appropriate statistical method. A simple example using the IPF (Iterative Proportional Fitting) algorithm [16] shall illustrate a possible weighting procedure, its successful application and possible failures.

Germany	Sum	Urban	Rural		Germany	Sum	Urban	Rural
Sum	401,823	258,919	142,904	← input	Sum	100.0%	64.4%	35.6%
Injured	397,671	257,694	139,977		Injured	99.0%	64.1%	34.8%
Fatalities	4,152	1,225	2,927		Fatalities	1.0%	0.3%	0.7%
<hr/>								
GIDAS	Sum	Urban	Rural		GIDAS	Sum	Urban	Rural
Sum	2,203	1,680	523		Sum	100.0%	76.3%	23.7%
Injured	2,178	1,676	502	← input	Injured	98.9%	76.1%	22.8%
Fatalities	25	4	21		Fatalities	1.1%	0.2%	1.0%
<hr/>								
<b>Iteration 1</b>	<b>Row</b>	<b>Col</b>						
	401,823.0	306,677.4	95,145.6	258,919.0	142,904.0			
	397,671.0	306,013.1	91,657.9	396,023.8	258,358.1	137,665.7		
	4,152.0	664.3	3,487.7	5,799.2	560.9	5,238.3		
<b>Iteration 2</b>	<b>Row</b>	<b>Col</b>						
	401,823.0	259,834.3	141,988.7	258,919.0	142,904.0			
	397,671.0	259,432.7	138,238.3	397,648.2	258,518.9	139,129.4		
	4,152.0	401.6	3,750.4	4,174.8	400.1	3,774.6		
<b>Iteration 3</b>	<b>Row</b>	<b>Col</b>						
	401,823.0	258,931.6	142,891.4	258,919.0	142,904.0			
	397,671.0	258,533.7	139,137.3	397,670.7	258,521.1	139,149.6		
	4,152.0	398.0	3,754.0	4,152.3	397.9	3,754.4		
<b>Iteration 4</b>	<b>Row</b>	<b>Col</b>						
	401,823.0	258,919.2	142,903.8	258,919.0	142,904.0			
	397,671.0	258,521.3	139,149.7	397,671.0	258,521.1	139,149.9		
	4,152.0	397.9	3,754.1	4,152.0	397.9	3,754.1		
<b>Iteration 5</b>	<b>Row</b>	<b>Col</b>						
	401,823.0	258,919.0	142,904.0	401,823.0	258,919.0	142,904.0		
	397,671.0	258,521.1	139,149.9	397,671.0	258,521.1	139,149.9		
	4,152.0	397.9	3,754.1	4,152.0	397.9	3,754.1		
<hr/>								
<b>Result of IPF: GIDAS sample adapted to German national statistics</b>								
Germany	Sum	Urban	Rural		Germany	Sum	Urban	Rural
Sum	100.0%	64.4%	35.6%		Sum	100.0%	64.4%	35.6%
Injured	99.0%	64.3%	34.6%		Injured	99.0%	64.3%	34.6%
Fatalities	1.0%	0.1%	0.9%		Fatalities	1.0%	0.1%	0.9%
<hr/>								
max error						0.2%		

Table 2 – IPF algorithm applied to GIDAS as a sample of German National Statistics.

As the name suggests, IPF adapts the frequencies (cell values) in a contingency table to a marginal distribution using an iterative process that must converge in order to get a result. Fortunately, in most real life scenarios, convergence of IPF is good and fast. Table 2 shows how IPF is applied for the two parameters “injury severity” (fatal, injured) and “location of accident” (rural, urban). Expanding this to n parameters leads to an n-dimensional crosstab that can be handled by IPF in the same way.

Accident data used for this example is a GIDAS in-depth sample and national statistics from Germany and Austria, all from the accident year 2009. Input is a crosstab “location of accident” vs. “injury severity” in GIDAS and four marginal distributions of each of the two parameters for Germany and

Austria. The example has some realistic aspects as the injury severity is often needed for assessment of the potential of a safety measure. Apart from the general availability of injury severity in the national statistics and thus its likely application in a weighting procedure, accurate estimates of the injury severity is desirable.

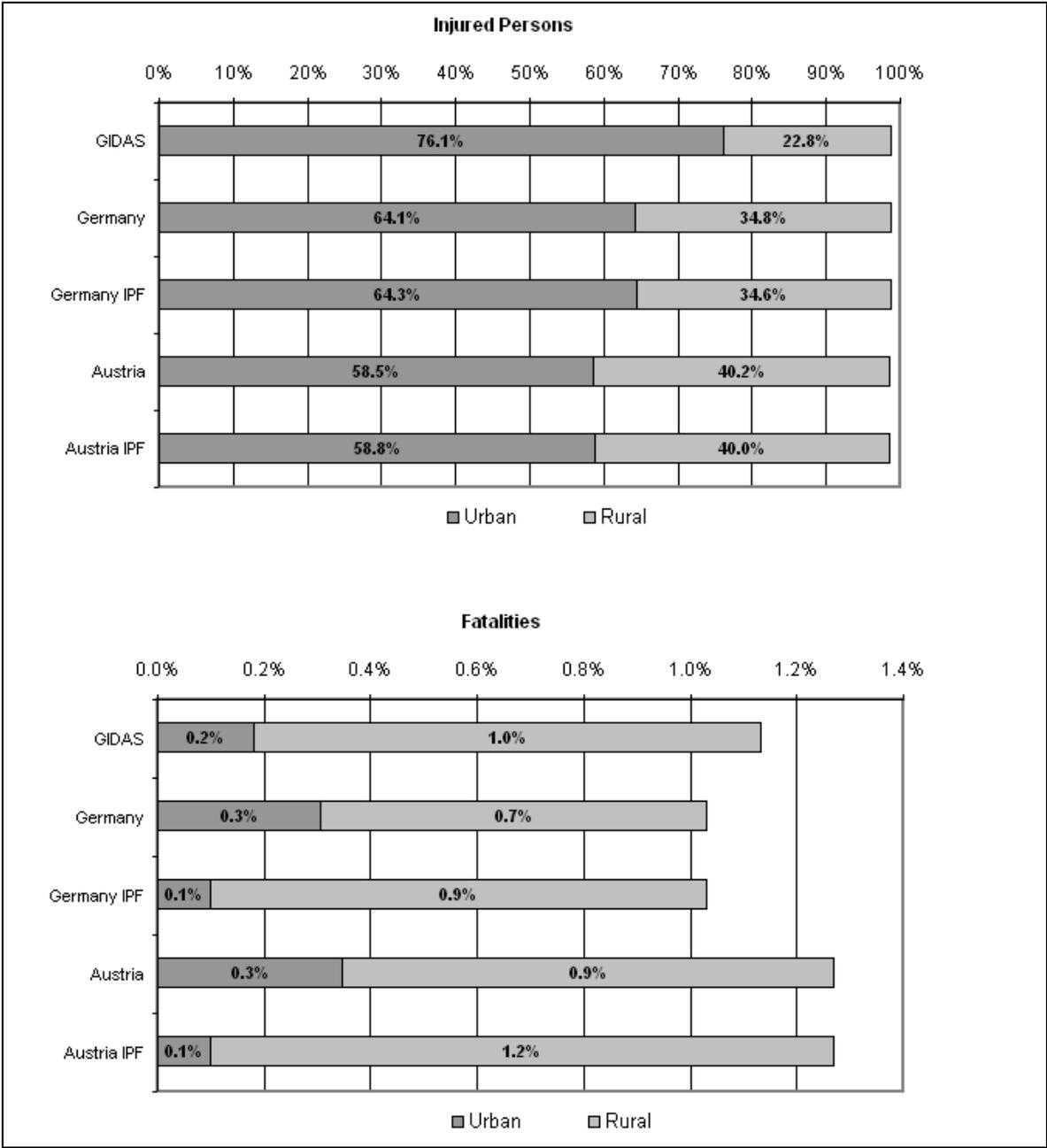


Chart 1 – Injury severity and location of accident in GIDAS sample adapted to Germany and Austria.

As the crosstabs for these two parameters are also provided for Germany and Austria, error checking can be conducted. Here, the overall error is calculated as the maximum relative error over all cells related to the whole sample size. Start of the iteration is a table combined of the marginal values of the national statistics and the frequencies of the combined parameter values in the GIDAS sample, which seed the starting table (top of table 2). Seeding has an influence on the cell values of the resulting crosstab, but not on its margins, they remain the same as in the starting table. Zero values in the starting table are a possible reason for non-convergence, as zeros are invariant throughout the iteration. This can be avoided by injecting small numbers replacing the zeros. In the Iteration step, alternating

rows and columns are calculated by weighting the corresponding value in the previous table with a quotient of the original marginal value and the one from the previous iteration step.

In this example, convergence is sufficiently reached after five steps. The resulting table is shown in the lower right corner of table 2 for the German national data and both results for Germany and Austria are shown in chart 1. Although there is bias in the sample, IPF provides a good estimate for the cells of the crosstab of the two weighting parameters. This is especially true for larger groups (injured persons). The maximum error rate amounts for 0.2% in the estimates for Germany and Austria, which is good in the context of the bigger group of injured persons. However, for smaller groups like the fatalities the error rate can have significant impact on the stability of the results, which then have to be interpreted carefully.

Finally, weighting factors can be derived as the quotient of percentages of the national estimates (after applying IPF) and the original in-depth sample percentages. Each accident in the sample then can be weighted with a factor that is given by the estimated crosstab cell value entry with the particular combination of weighting parameter values for this accident. This step is not very revealing in this small example, as more than two weighting parameters should be involved for acceptable results.

Of course, the more interesting situation where parameters not contained in the national statistics are estimated (a typical example may be airbag deployment) is likely to introduce additional errors with the constraint that these errors cannot be determined, even in a test case.

A reason for the results being quite good in this example is the inner relation between “location of accident” and “injury severity”, which is usually also not affected by regional differences. Reasons for this relation are higher speeds and more single car accidents in rural areas.

Adapting to a sample is theoretically not a matter of national borders. Hence, work is also underway to provide cluster methods to group similar countries in terms of accident data characteristics. This also has potential to increase the size of the in-depth subsamples.

## 5. NEXT STEPS, LONG TERM GOALS

A second working group meeting is planned for September 2012. There, first results of the pilot study will be presented and discussed and steps towards a finalization of the common data scheme will be taken. Having accomplished this, the next step is to build data converters, to transform the data into the common data subset. In some cases this may involve some contract work to be accomplished, leaving funding issues to be addressed. However, most of the work should be done self contained by the iGLAD working group and its associated organizations, reflecting the light-weight nature of the project. IGLAD’s long term goal is the official standardization of the data scheme and a certification procedure addressing quality maintenance issues of the data set.

## 6. REFERENCES

- [1] German In-Depth Accident Study, <http://www.gidas.org>
- [2] EU Commission, European Road Safety Action Programme, ISBN 92-894-5893-3, 2003.
- [3] WHO, [http://www.who.int/roadsafety/decade\\_of\\_action/en/](http://www.who.int/roadsafety/decade_of_action/en/)
- [4] D. Ockel, J. Bakker, R. Schöneburg, “Internationale Harmonisierung von Unfalldaten; Fortschrittsbericht des FIA / ACEA Projekts iGLAD (Initiative for the Global Harmonization of Accident Data)“, VDI Kongress 2011, Berlin.
- [5] STAIRS, EC 4<sup>th</sup> Framework Programme, “Standardisation of Accident and Injury Registration Systems”, Final Report, Contract N° : RO-96-SC.204,1999
- [6] EACS, ACEA / EC project 1996 – 2001 coordinated by CEESAR, “European Accident Causation Survey - The Databank Questionnaire”, 2001
- [7] Pendant, <http://www.vsi.tugraz.at/pendant/>
- [8] TRACE Project, <http://www.trace-project.org>

- [9] SafetyNet, <http://erso.swov.nl/safetynet/content/safetynet.htm>
- [10] DaCoTa, <http://www.dacota-project.eu/>
- [11] MAIDS, <http://www.maids-study.eu/>
- [12] ETAC, EU Commission, “European Truck Accident Causation – Final Report”, 2006
- [13] IRTAD, <http://internationaltransportforum.org/irtadpublic/index.html>
- [14] CARE, [http://ec.europa.eu/transport/road\\_safety/specialist/statistics/index\\_en.htm](http://ec.europa.eu/transport/road_safety/specialist/statistics/index_en.htm)
- [15] CADaS, SafetyNet Integrated Project D.1.14, “The Common Accident Data Set”, 2008
- [16] W.E. Deming, F.F. Stephan, On a Least Squares Adjustment of a Sampled Frequency Table When the Expected Marginal Totals are Known, *Annals of Mathematical Statistics* 11 (4): 427–444. doi:10.1214/aoms/1177731829. MR3527, 1940