

# **Estimating Data Volumes of RFID-enabled Supply Chains**

**Alexander Ilic**

Information Management  
ETH Zurich, CH-8092 Zurich  
ailic@ethz.ch

**Andrea Grössbauer**

Information Management  
ETH Zurich, CH-8092 Zurich  
andreagr@student.ethz.ch

**Florian Michahelles**

Information Management  
ETH Zurich, CH-8092 Zurich  
fmichahelles@ethz.ch

**Elgar Fleisch**

Information Management,  
ETH Zurich, CH-8092 Zurich, and  
Institute of Technology Management,  
University of St. Gallen, CH-9000 St. Gallen  
efleisch@ethz.ch

## **ABSTRACT**

The widespread application of RFID tags in supply chains is said to cause enormous data volume problems and thus unprecedented challenges for systems and infrastructures. In order to unleash the potential of item-level RFID applications, such as data sharing and discovery across company boundaries, an unbiased understanding of emerging data volumes is necessary. However, quantitative data that provides factual argument is still scarce. Therefore, we present a simulation study based on a real-world scenario that reveals quantitative characteristics of the data volumes problem in an RFID-enabled supply chain and discuss its implications. Our results suggest that data volumes will be much lower than currently anticipated, but still bear significant challenges for researchers and developers of RFID infrastructures.

## **Keywords**

Radio Frequency Identification, data volumes, supply chain, EPC network, simulation.

## **INTRODUCTION**

Radio Frequency Identification (RFID) technology in combination with standardized numbering schemes such as the Electronic Product Code (EPC) enables the tracking of physical objects in supply chains. Raw ID reads of RFID tags are enriched by context (time, reader location, etc.) and captured as so-called RFID events. These RFID events are a key enabler for increasing efficiency of various business processes and allow for a more fine-granular way of management. Retail giants (e.g. Wal-Mart), aerospace companies (e.g. Boeing), pharmaceutical companies and many other companies are interested in leveraging these benefits and are involved in creating a global set of standards for RFID-based business applications, the so-called EPC Network (EPCglobal, 2007a). One of the critical parts of this set of standards is the EPC Information Services (EPCIS) specification for the management of RFID events (EPCglobal, 2007b). Although there is common agreement that sharing RFID event data could lead to significant efficiency gains, often terms like “enormous data volumes” (Vadlamani, Kalyan and Murali, 2006), claims of “7 TB” (Schuman, 2004) or “15 TB” (Bhuptani and Moradpour, 2005) per day generated by a large retailer alone, or over-simplified calculations (e.g. each item leaves 10 traces per day (Gonzalez, Jiawei, Xiaolei and Klabjan, 2006) are used to indicate that current technology is by far not capable of coping with these challenges. Up to our best knowledge, no scholarly paper exists that provides quantitative and reliable data for supporting these claims and analyzing the impacts on individual supply chain members. In fact, many organizations claim that the impact of expected data volumes are a major source of uncertainty that hinders RFID adoption. In addition justified figures about expectable data volumes could support researchers and developers of RFID event databases, discovery services, access control frameworks, and analytical tools in their design decisions and performance assessments.

Thus, the goal of this paper is to provide a quantitative data volume analysis for RFID-enabled supply chains. We estimate data volumes at the example of a retail supply chain since the retail industry is considered being a lead adaptor of RFID technology. As the generation of RFID events is closely linked to characteristics of product flows and replenishment management in a supply chain, we apply a simulation modeling methodology based on parameters of a real-world supply chain.

The paper is structured as follows. In the next section, we introduce the implemented supply chain model and its assumptions. Afterwards, we present and explain the results from the simulation runs. Then, we extend the results to multiple products and stores and thereby assess the impacts on a complete retail supply chain. Then, we discuss our findings and describe the most important lessons learned. The final section concludes this paper and highlights future work.

## SIMULATION MODEL

### Business scenario and assumptions

The simulation model (Figure 1) bases on the characteristics of a real-world scenario and represents a typical retail supply chain in Europe. Simulation parameters were derived from national sales and shipment figures from the years 2006 and 2007 for a selected food product in the retail industry and interviews with logistics managers, CIOs, and supply chain experts of the involved parties.

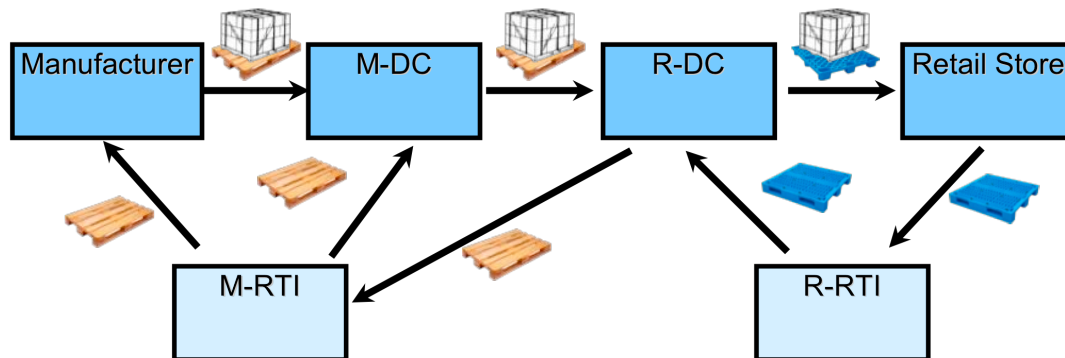


Figure 1. Flow of goods and pallets in the simulation model

The product is a typical food product with a shelf-life of one year. Thus, the aspect of perishability can be neglected in our simulation. The demand rate (units sold per day at the retail store) or short “demand” is Poisson distributed (Zipkin, 2000), averages 5.5 items per day in the base case, and is subject to seasonality factors. The seasonality pattern was extracted from actual sales data according to the method of Chopra and Meindl (2007) and comprises eight seasonal periods. Demand is satisfied according to a service in random order (SIRO) policy. Lost sales occur, if due to empty stock at the retail store consumer demand cannot be fulfilled (Zipkin, 2000).

The flow of goods is depicted on Figure 1. A retail store sells goods to consumers and receives replenishments from the retail distribution center (R-DC) once the stock level is below a certain level. The inventory of the R-DC is managed by the vendor. In this vendor managed inventory approach (Ghiani, Laporte and Musmanno, 2004), the manufacturer replenishes the retail distribution center with shipments from the manufacturer’s distribution center (M-DC). The M-DC receives shipments from the manufacturer. In order to optimize replenishment decisions and minimize out-of-stock situations, the R-DC is replenished according to an adaptive forecasting technique based on historical sales data (Chopra and Meindl, 2007).

For each delivery step, we assume positive, normally distributed lead-times. The lead-times in Table 1 base on settings of our case study partners and represent typical values for an European supply chain. Moreover, pallets are rented from the manufacturer’s pallet pool provider (M-RTI) and the retailer’s pallet pool provider (R-RTI) in order to transport the goods in the delivery steps. In addition, an internal manufacturing lead-time of 6.00h (stdev: 1.2h) occurs to account for set-up times in the production line. We assume that the size of the pallet pool of M-RTI and R-RTI is unlimited and thus pallets are always available at the required sites. As we investigate the product flow of a single product class, we assume homogenous

pallets only. Once a number of 40 empty pallets accumulates at a site, they are returned to the corresponding pallet pool providers (pallet collection process). Figure 1 shows the flow of goods and the supply and return of pallets.

Source	Destination	Lead-time (mean)	Lead-time (stdev)
Manufacturer	M-DC	1.33h	0.17h
M-DC	R-DC	8.00h	0.33h
R-DC	R	0.67h	0.08h

Table 1. Delivery lead-times in simulation model based on case study data

### RFID set-up and event data model

For our simulation, we consider the following three tagging levels based on EPC numbering schemes of the GS1 traceability standard (GS1, 2006):

- Pallet level: In pallet level tagging, a RFID tag is integrated into the pallet and identifies it via a Global Reusable Asset Identifier (GRAI). Moreover, in order to track the association between load and pallet, the logistic unit, for a shipment, a Serial Shipping Container Code (SSCC) is used in conjunction with an additional RFID tag attached to the shrink-wrapped packaging. Alternatively, this identifier can be programmed into the memory of the pallet's tag.
- Case level: In addition to pallet level tagging, an RFID tag is applied to every trade unit. Each case is identified by a Serialized Global Trade Identification Number (SGTIN).
- Item level: In addition to case level tagging, an RFID tag is applied to every single consumer unit, which is also identified by a Serialized Global Trade Identification Number (SGTIN).

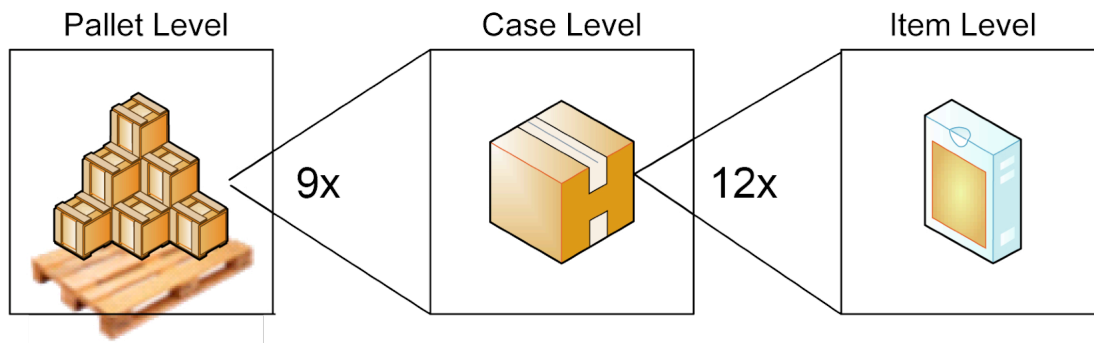


Figure 2. Packaging hierarchy of the case study with corresponding cardinalities

The hierarchical relationship between pallets, cases, and items is depicted on Figure 2. Depending on the selected tagging granularity a RFID reader has to read several tags and create corresponding events. We assume that at all six supply chain sites of Figure 1, sufficient RFID readers and read points are deployed in order to generate these RFID events. Recorded events follow the EPCIS event standard specification (EPCglobal, 2007b). EPCIS events can capture properties of an RFID read with respect to the dimensions of “what?”, “when?”, “where?”, and “why?” and typically represent them in XML format. Events for different items showing the same characteristics (e.g. timestamp, location, business step, ...) are grouped into a single events according to typical filtering and aggregation practices (EPCglobal, 2007a). In the simulation, we use the three different EPCIS event types described in Table 2.

EPCIS event type	Description	Event size
ObjectEvent	Captures contextual information pertaining to one or more physical objects identified by EPC codes.	>84 bytes
AggregationEvent	Describes the physical aggregation or disaggregation of objects which are identified by EPC codes. Relations between parent and child objects are transitive.	>84 bytes
TransactionEvent	Relates to the EPCs involved in the same business transaction. Unlike aggregation events, transaction events describe “weak” associations (e.g. multiple pallets being part of the same shipment)	>104 bytes

Table 2. Properties of event types used in simulation and minimum event sizes

## SIMULATION RESULTS FOR A SINGLE PRODUCT SCENARIO

### Setup

The simulation model was implemented in Python using the SimPy discrete event simulation framework. The design of experiments comprised running the previously described base scenario for a simulated time of 595 days in 10 independent replications. Because the simulation is initialized with empty stock levels, a warm-up time is needed in order to reach a reliable state. The length of the warm-up period was set to 230 days according to the method described in (Law and Kelton, 2001) in order to fulfill this condition. For the remaining 365 days, the performance measures of number of events and data volumes generated at each supply chain location were observed. In order to assess the sensitivity of the results, additional simulations with different tagging levels and demand levels were conducted.

### Results

Figure 3 shows the results of the simulation runs with different tagging levels in a single product scenario. For each supply chain party, the observed data volume is shown dependent on the simulated tagging granularity. The impact for the pallet pool providers M-RTI and R-RTI is independent from the tagging granularity. The reason is that only the tagging granularity was varied and other relevant parameters, such as consumer demand, were left unchanged. As a result, the same number of outgoing and incoming pallets is needed. The small deviations in the graph are caused by variations in the random number generator. For the other parties, there is a significant increase in data volumes when switching from pallet to case level tagging. The reason is that additional data is needed in order to track the association between pallets and cases. When moving to the item level, there is, however, no additional data increase for M-DC and R-DC. The reason is that the physical association between items and cases is established at the manufacturing level and not revoked until reaching the retail store (i.e. trade units are opened) and thus no additional events are required.

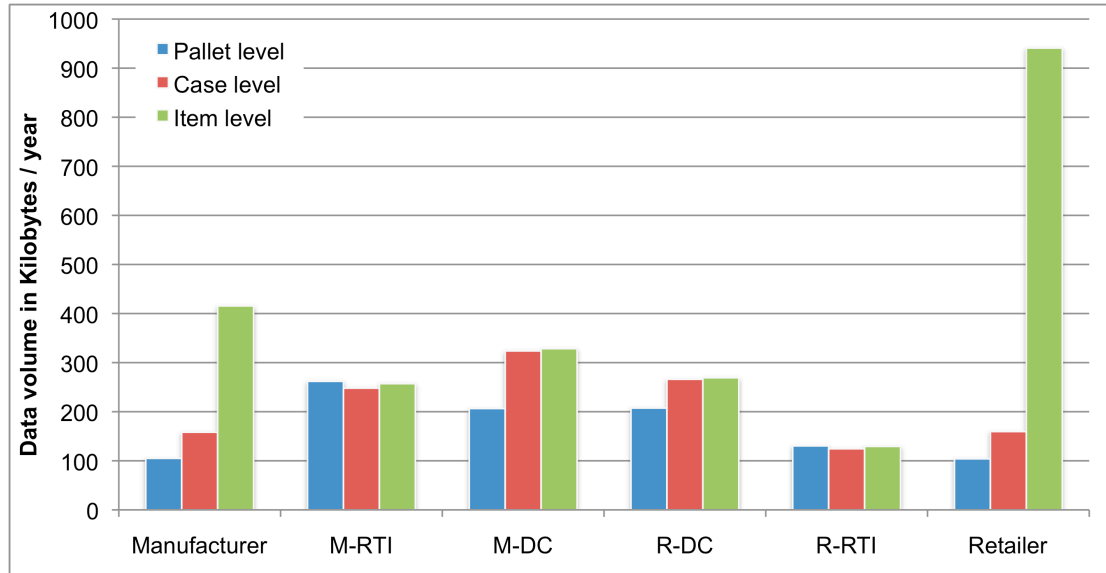


Figure 3. Simulation results showing data volumes generated by each partner for different tagging granularities

As explained above, a switch to a more fine granular tagging level leads to a shift in the data volume distribution. A surprising result is that the strongest impact of moving from case to item level will be on the manufacturer and the retailer, as every individual item needs to be read (see Figure 4). The high peak at the retailer compared to the manufacturer is not intuitive at first glance, but can be explained quite easily. Based on the principle of basic aggregation and filtering, the manufacturer can group events because items are manufactured in batches. In contrast, the retailer sells single items dependent on customer arrival. Thus events cannot be grouped anymore as, for instance, their timestamp is completely different.

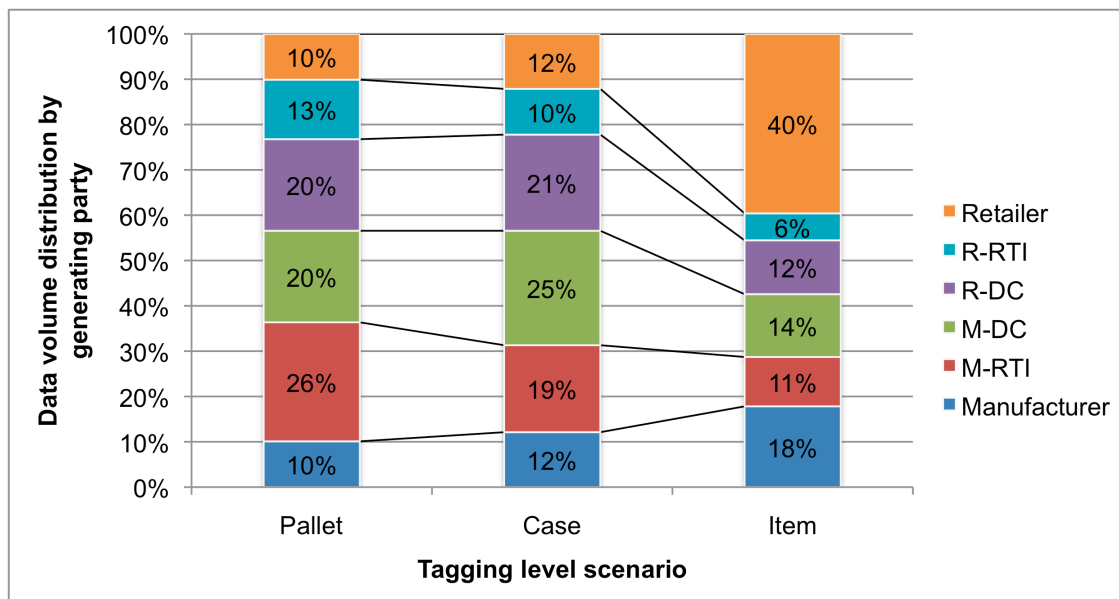


Figure 4. Shift in data volume percentages by tagging granularity

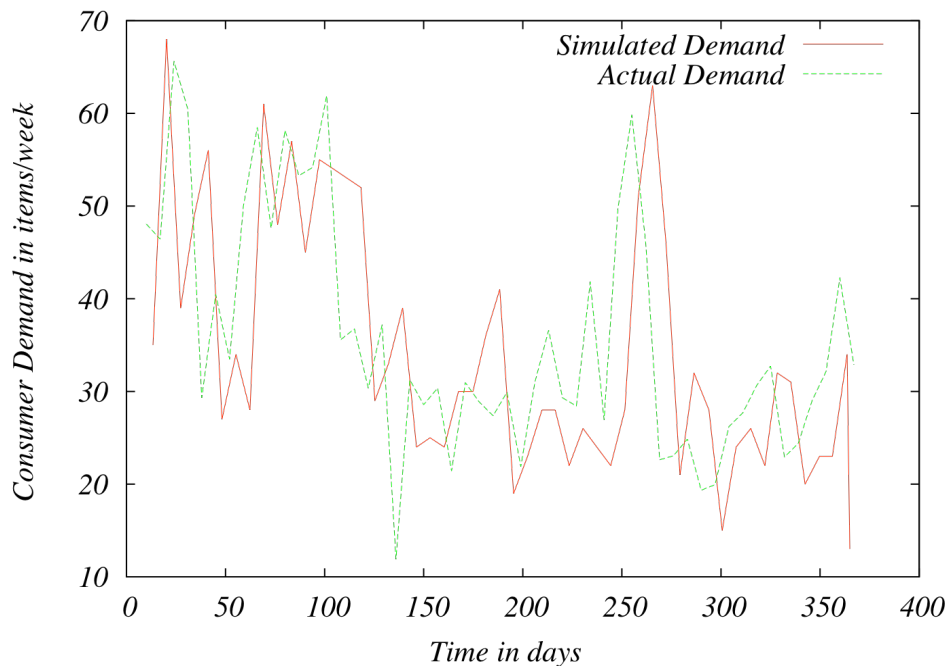
## Verification and Validation

Regarding verification (Sargent, 2005) of the model, we created flow charts of the individual process steps and data models in order to verify the implemented parameters, logic, and product flows with supply chain experts. In order to ensure reproducible and stable results, we applied a common random number (CRN) approach (Law et al., 2001) and executed 10 replications per individual simulation run.

To validate our simulation model, we implemented the supply chain characteristics in a computer model and then validated the results with different techniques. The validation tests conducted are depicted on Table 3. Details how to conduct the corresponding tests are described well by Sargent (2005). For each of the tests, we compared historical data of the specific product with the simulation results and asked an expert whether the simulation model inhibits the same relevant characteristics. As an example, Figure 5 shows the output of simulated demand at the retail store versus the historical data. Together with a supply chain expert we could confirm that this operational graphic sufficiently mimics the demand patterns and its seasonally induced fluctuations (Ghiani et al., 2004).

Item	Pallet flow	Customer demand	Seasonal pattern	Forecasting process
<b>Validation technique</b>	Face validity	Face validity	Operational graphic	Face validity
<b>Reference data</b>	Shipment data	Sales figures	Sales figures	Order history

**Table 3. Validation techniques and reference data used**



**Figure 5. Simulated demand versus historical data**

## Demand sensitivity

The previously described simulation model focuses on a single product scenario and thus shows relatively low data volumes (see Figure 3). In reality, a manufacturer has multiple products, each with different demand patterns. As the demand pattern for e.g. a can of soda could be completely different from milk, it is important to assess the sensitivity for varied demand levels to estimate the data volumes of multiple product classes. Figure 6 shows the total data volume when varying demand

for the three different tagging levels. The results suggest that there is a linear relationship between the mean demand level and the total data volume. This allows us to extend the findings of Figure 4 to a multiple product scenario under the assumption that products differ mainly in their demand characteristics.

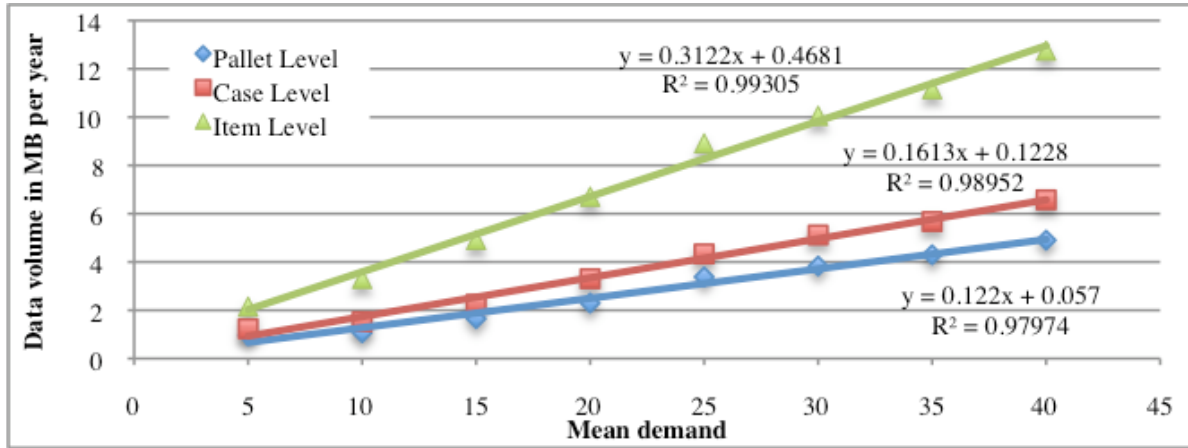


Figure 6. Total supply chain data volume dependent on tagging level and demand

#### EXTENSION TO A MULTIPLE PRODUCT SCENARIO

In this section, we extend the scope from a single to a multiple product scenario. This allows us make a reasonable impact assessment on a complete retail supply chain. The impact analysis focuses on item-level tagging, since it is said to offer tremendous efficiency gains for retailers and real-world deployments in the retail industry are still scarce.

##### Data volumes by retail store size

Typically, one can distinguish between three different retail store sizes: small, medium, and large. The main difference between the store types is the number of sold items per day. An expert interview revealed that for a small store 1'800, for a medium store 8'000, and for a large store 30'000 consumer units are typically sold per day.

Assuming that the simulated product represents an average retail product, we can extend the impact assessment to multiple products and thus assess the impact on the previously described store sizes by taking the number of consumer units sold per day as mean demand parameter of our simulation. The results are shown on Figure 7. The data volume generated through the replenishment of a large retail store amounts to 9.15 GB per year. Over 40% of the data are generated at the retail store, while the remaining 60% are generated throughout the supply chain by other parties. The IT system of a large retail store must be able to handle on average 12 MB per day, which can be easily achieved with current technology.

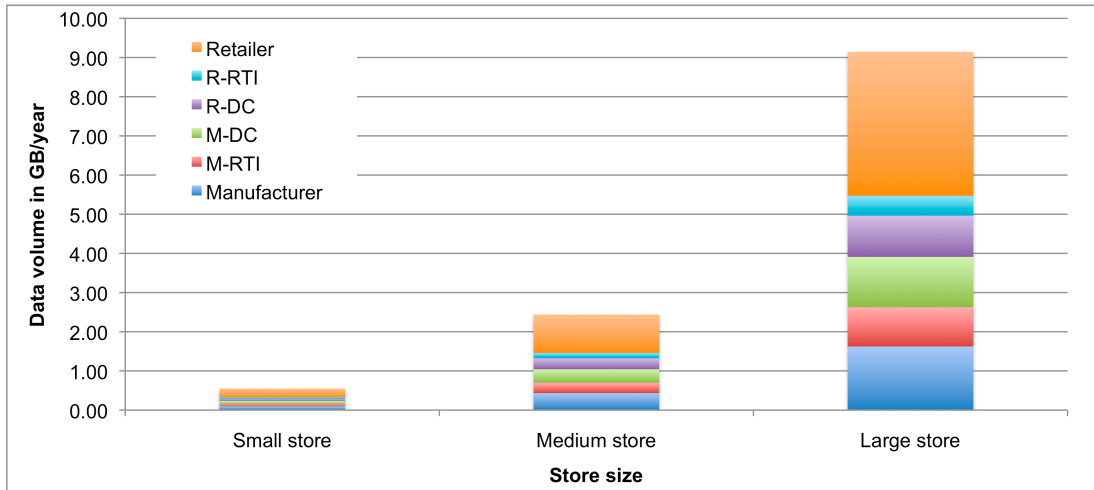


Figure 7. Data volumes in the supply chain for replenishing different retail store sizes

#### Total data volume generated by a retail supply chain

A retail supply chain has a number of distributed stores with different sizes. However, event data is usually collected from stores and supply chain partners, and then analyzed in a central IT system. In the following, we take the example of a retailer in Switzerland with the corresponding number of different stores. Table 4 shows that per year, the central IT system of the selected retailer has to handle 1042 GB of RFID event data.

Store size	Number of stores	Data volume per store/year	Total data volume/year
Small	290	0.55 GB	159 GB
Medium	212	2.44 GB	517 GB
Large	40	9.15 GB	366 GB
<i>Total</i>			1042 GB

Table 4. Total data volume calculation from an example retailer's perspective

#### DISCUSSION AND LESSONS LEARNED

This section aims to capture the most important lessons learned during the iterative process of refining the simulation model and the discussions of our findings with industry experts.

##### Implications for sizing the problem space

Our exemplary data volume estimation for a full retail supply chain revealed that item-level tagging would generate 3.3 GB per day (1042 GB per year) of event data. When comparing the figure of our example supply chain with the often cited supply chain of Wal-Mart with its over 2580 Supercenters, the frightening number of 7 TB per day (Schuman, 2004) seems to be largely overestimated. In fact, when comparing the 3.3 GB per day of the previous' section analysis to the retailer's point-of-sale (POS) data of 2007 (approximately 1.8 GB per day), it becomes obvious that the analysis of item-level RFID data could be done with a similar sized infrastructure. However, unlike POS data, which is often transmitted and evaluated in batches (e.g. daily), the real-time decision-making vision of RFID could still make it challenging to handle the data volumes in an efficient way. The challenge gets even bigger when considering data overheads and increased complexity caused by the fact that RFID data is often spread over multiple locations, supply chain partners, and systems.



### **Implications for research and standards development**

Current RFID standards on information sharing do not yet assign a high priority on data volumes. However, a sound approach is needed in order to make proper architectural decisions and to understand the nature of full-scale item-level tagging in detail. So far, our study was limited to the retail supply chains. Early talks with researchers and experts of other industries already provided some hints that data volumes and requirements might be different in other industries and thus demand for a more targeted approach.

### **Relevance of adoption tools and simulation**

Item-level tagging is still too costly to roll out in retail supply chains. Our simulation model allowed assessing the impact of a switch from pallet to item level tagging in a cost efficient way while retaining the business characteristics of the scenario. The advantage of a simulation approach is that all parameters can easily be adapted to other scenarios and thus generate specific learnings for each individual case. In order to allow companies to test their systems and information sharing strategies for the future item-level scenarios, we implemented our simulation model so that it generates valid EPCIS data in XML representation. With this adoption tool, companies can benchmark their systems and infrastructures and also vary changes in their specific setting. We aim to make the simulation tool publicly.

### **Relevance of information-sharing architectural design**

During the design of our simulation model several discussions emerged about the information sharing infrastructure. So far, the simulation model of above just counts data generated at the different supply chain locations. In reality, however this data must somewhere be consolidated and made available for querying (e.g. for product recall purposes etc.). We deliberately did not consider the data usage aspect in our simulations since multiple suggestions for an infrastructure design exist (see e.g. (Beier, Grandison, Kailing and Rantzau, 2006), (Kurschner, Condea, Kasten and Thiesse, 2008) for more details). Future research will have to investigate the aspects of querying, indexing, processing with respect also to data throughput and data storage implications.

### **Relevance of filtering and data compression**

There are several levels of filtering in RFID infrastructures. In our simulation, we assumed that the huge number of low-level and middleware events are already filtered and aggregated into fewer, higher-level EPCIS events. A deliberate information sharing strategy should, however, ensure that only business relevant events are stored in an EPCIS information sharing repository and thus reduce the number of events already at the source. Since most EPCIS standard implementations rely on the XML binding, EPCIS events are often stored in an inefficient way and contain much redundant information. For our simulation, we just investigated the XML data volume, which is generated at each partner's facility. Compression tests of the simulation output XML showed that the size of an EPCIS repository could easily be reduced by a compression factor of 45 and thus providing significant leverage due to intelligent data storage (e.g. (Gonzalez et al., 2006)) and transmission techniques.

## **CONCLUSIONS AND FUTURE WORK**

The goal of this paper was to provide a quantitative analysis of RFID event data volumes in a supply chain. By using a simulation modeling approach based on characteristics of a real world supply chain, we made a first step in this direction and were able to start from a proven basis toward future scenarios such as case and item-level tagging without needing to conduct costly field trials. The positive responses from our industry partners confirmed the value of this approach and helped to derive lessons learned for an envisioned field trial. As our approach bases on the characteristics of an average, single type of product, the findings must be interpreted carefully for scenarios with completely different characteristics. However, our results show that the data volume problem induced by item-level tagging is likely to be overestimated, but still bears significant challenges for processing and analysis of RFID event data in supply chains. One of the key learnings is that for a retailer, the data volumes are likely to be in the same magnitude of order as POS data is today. Thus, with hardware of today it should be possible to cope with data volumes of item-level tagging. However, efficient inter-organizational information sharing still comprises several significant unresolved software challenges due to real-time processing and distributed data management requirements. Researchers and developers of information sharing repositories, data discovery services, and event data storage systems can potentially use our figures and simulation approaches as a first starting point for optimizing design decisions and conducting performance evaluations.

Future research on RFID event data volumes is still needed in order to verify and extend our results to different types of products and industries. Once data from field trials on case and item-level tagging is available, simulations combined with other approaches such as value stream mapping can help to optimize data flows even further by identifying bottlenecks and critical parts. Different designs and architectures for core information sharing components such as data repositories or discovery services can then be evaluated on a more quantitative basis. Since the number of scholarly articles in this area is very limited, there is significant potential for discovering new insights and implications that can clearly shape and direct the development of future RFID infrastructures.

## REFERENCES

1. Beier, S., Grandison, T., Kailing, K., and Rantzau, R. "Discovery Services—Enabling RFID Traceability in EPCglobal Networks," 2006.
2. Bhuptani, M., and Moradpour, S. "Emerging Trends in RFID," Jupitermedia, Darien, USA.
3. Chopra, S., and Meindl, P. *Supply Chain Management: Strategy, Planning, and Operation*, (3rd ed.) Prentice Hall, Upper Saddle River, New Jersey.
4. EPCglobal "Architecture Framework v. 1.2," EPCglobal, Brussels, Belgium.
5. EPCglobal "EPCIS Standard v. 1.0.1."
6. Ghiani, G., Laporte, G., and Musmanno, R. *Introduction to Logistics Systems Planning and Control* Wiley, 2004.
7. Gonzalez, H., Jiawei, H., Xiaolei, L., and Klabjan, D. "Warehousing and Analyzing Massive RFID Data Sets," in: *Proceedings of the 22nd International Conference on Data Engineering*, 2006, pp. 83-83.
8. GS1 "The Global Traceability Standard," GS1, Brussels, Belgium.
9. Kurschner, C., Condea, C., Kasten, O., and Thiesse, F. "Discovery Service Design in the EPCglobal Network - Towards Full Supply Chain Visibility," in: *The Internet of Things Conference*, Zurich, 2008.
10. Law, A.M., and Kelton, D.M. *Simulation Modeling and Analysis* McGraw-Hill Higher Education, 2001.
11. Sargent, R.G. "Verification and validation of simulation models," in: *Proceedings of the 37th conference on Winter simulation*, Winter Simulation Conference, Orlando, Florida, 2005.
12. Schuman, E. "Will Users Get Buried Under RFID Data?," Ziff Davis Media, New York, USA.
13. Vadlamani, R., Kalyan, R., and Murali, K.M. "Applications of Intelligent Technologies in Retail Marketing," in: *Retailing in the 21st Century*, 2006, pp. 127-141.
14. Zipkin, P.H. *Foundations of inventory management* McGraw-Hill, Boston, 2000.