

# Procurement data in context: Analysis of the Procurement Process in the area of coatings and paints <sup>\*</sup>

Business Process Intelligence Challenge 2019

Urszula Jessen<sup>1</sup>

akquinet cubit GmbH [urszula.jessen@akquinet.de](mailto:urszula.jessen@akquinet.de)

**Abstract.** Understanding and managing the processes in a sustainable and optimized way is one of the significant challenges in modern organizations. Process mining offers both methodologies and tools to get unique insights into the process data. The process discovery, visualization, and conformance checking provide support not only for the management or process owners but also for all process users. On the other hand, most companies already possess different tools and paramount knowledge on how to analyze and understand data in their specific domain. In our approach, we try to find synergies between the classic data and process mining to let the users with tools and domain knowledge to deploy it in their daily tasks. In the following paper, we will demonstrate the framework for combining Process Mining Techniques with Data Science and BI Tools to achieve the best data insights for the given context and to utilize existing investment in knowledge and software in the organization. Based on the data from a large multinational company operating in the area of coatings and paintings, we develop a set of tools for exploratory analysis of the purchasing process.

**Keywords:** BPI Challenge · Process Mining · Event log · Process Intelligence · Descriptive Analytics.

## 1 Introduction

Business Process Intelligence Challenge is a yearly competition for process mining students, researchers, and professionals. The data provided in BPI Challenge 2019 originates from a large multinational company operating from the Netherlands in the area of coatings and paints and covers parts of its purchase order handling process. Each Process Order in Event log contains one or multiple line items.[1]. This data has been collected for the purchase order handling process for some of companies 60 subsidiaries. The process owner has some questions regarding the process and is generally concerned about compliance issues. The other problems that need to be addressed are generalized process models, process efficiency, process outliers and deviation from the standard process flow.

---

<sup>\*</sup> Supported by Process Science.

The starting point of our analysis is understanding the processes and the context for the provided data. We will split the data into separate categories that can be described in a general manner as similar processes. We will present some of the "birds-eye" process maps with its most essential characteristics. After establishing the general framework for further research we will focus on three aspects of purchase processes: Efficiency, Effectiveness and Compliance. We will then design measurements and performance indicators for those aspects of the process. Among others, the time, costs and complexity of the process will be examined. We will separately address the questions of the process owner. In the end, we present our conclusions and further recommendation for this project.

Although we gave our best efforts to understand main issues and context of the provided data, we are neither the experts in the given domain nor did we have chance to conduct interviews with the process owner and process users. For the analysis, we made general assumptions, so that some of the outcomes may differ according to individual company policies or rules. The challenge encouraged participants to use different tools, techniques, and methods. To analyze the data, we used various commercial and open-source tools. Apart from process mining techniques, we utilized advanced BI Tools and methods. The process owner has some questions regarding the process and is particularly concerned about compliance issues. The other problems that need to be addressed are generalized process models, process efficiency, process outliers, and deviation from the standard process flow.

## 2 Understanding provided data

The data presented gives an overview of the transactional purchase process, which covers the activities from creating purchase order till invoice clearance.

Although, the data should contain only the purchase orders from 2018, individual values in timestamps of the event log stretch out to 1948 in the past and 2020 in the future. As we cannot check if such values exist due to conversion error or if there are some exceptional cases in the system, we have decided to filter the timestamps to values from 1/1/2017 till 25/5/2019. We did not filter out the whole traces, but only the timestamps that did not match in the proposed timeframe. The resulted event log contains more than 99% of the original events, and all the traces from original data are preserved. In further analysis, all statements are based on this filtered data.

The provided data consists in total of 251.734 cases. Most of them contain Vendor invoice message (209.889), and goods receipt message (234. 479). Only 73% of the invoices are cleared, and more than half of all process flows contains repetitive tasks. The other characteristics of the process are a relatively low proportion of the Catalogue based purchases; only one per cent of instances include SRM related tasks. The throughput of the process flow is, on average, 64 days.

## 2.1 Classification of data

In the provided data process owner has specified four different process flows, 3-way matching, invoice after goods receipt, 3-way matching, invoice before goods receipt, 2-way matching (no goods receipt needed) and consignment. In 3 way matching an invoice is usually matched to the corresponding purchase document for quantity and value. Table 1 gives an overview of case statistics for the different flow categories. For better clarity, we applied the abbreviations for each of the flow categories.

**Table 1.** Different process flow types.

Case Item Category	Abbreviation	# Occurences and (% Occurences)	Description
3-way matching, invoice after goods receipt	3WIAGR	15.148 (6%)	In this type of process flow, the company expects matching values for the goods receipt message, an invoice receipt message and the value that was defined in the first item creation .
3-way matching, invoice before goods receipt	3WIBGR	220.814 (88%)	A goods receipt message is required, but there is no need for GR-based invoicing.
2-way matching	2W	891 (0.4%)	For these items, the value of the invoice should match the value at creation, but there is no separate goods receipt message required.
Consignment	CONS	14.494 (6%)	For these items, there are no invoices on PO level as this is handled fully in a separate process.

Apart from that event log contains multiple attributes like document type, spend area text, source, or spend classification. The purchase documents can also be divided into product related and non-product related goods and services. The NPR, also known as indirect goods and services include all goods and services that are not directly involved in the production, such as capital equipment, marketing, legal assistance or telecommunication. The NPR related process variants should be considered separately as they usually have some different characteristics from typical product-related processes [2]. This kind of purchase is generally time-consuming as the items are typically non-standardized and purchased in small bulks. The NPR usually needs much more user attention and not only from the purchasing department but also all different parts of the company. Figure 1 shows the connection between spend classification, spend are and flow type. In the provided data almost 70% of all orders belongs to product-related category.

The other feature dividing the purchase document can case Document Type, where there are three types of processes, Framework order, EC purchase order, and Standard PO. Especially Framework orders which are processes that have some specific validation time (with a start- and end- date) and item limits, differ in their characteristics from standard PO. This kind of purchasing process is usually used for services, travel expenses or utilities.

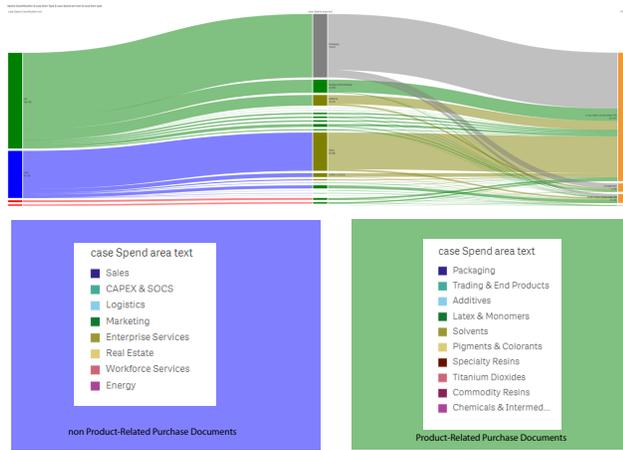


Fig. 1. Spend classification connection with spend areas and process flow types.

## 2.2 Data characteristics

Besides the typical data classification, it is essential to understand what kind of attributes should be considered as critical indicators for further inspection. In our evaluation, we want to focus on three aspects of the process: efficiency, effectiveness and compliance. Considering those aspects and available data, we concentrate on two distinct indicators: purchase volume and process throughput.

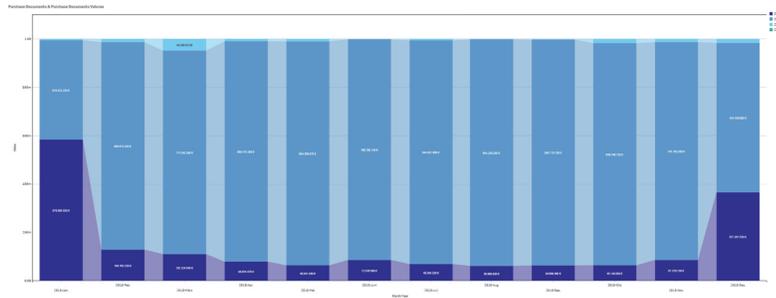


Fig. 2. Purchase volume for different type categories.

Figure 2 depicts the distribution of purchase volume throughout the year and for the different category types. For this analysis we used the average value of all documents in each process instance. Although the documents of the kind “3-way match, invoice before GR” have the most purchase documents line items and the most significant amount, there are two distinct value peaks for the document of

the type “3-way match, invoice after GR” in January and December. It can be explained through service or line items, that are usually paid yearly (at the beginning or the end of each year). Similar amount peak can be observed for “2-way match” documents in March. Overall is the purchase value stable and constant over the year. Consignment documents are not shown as they have no amount value.

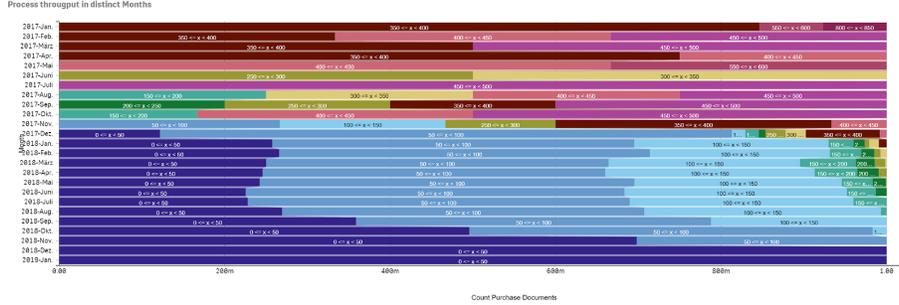


Fig. 3. Throughput of process instances(by year/month).

In the following analysis, we have divided throughput times into buckets of 50 days. Based on this classification, figure 3 depicts the throughput times of process instances from January 2017 till January 2019. There can be a distinct trend observed, in which process flow is getting steadily lower. Whereas most of the process instances in January 2017 have cycle times of 350-400 days, in January 2019 most of the duration of the case is not higher than 50 days. The most significant difference can be observed from January 2018.

### 3 Analyzing control flow

Process mapping is one of the most important parts of process mining projects. The literature describes multiple algorithms and tools for this task. In our analysis, we wanted to find the most generalized model of the given data and figure out how well does it fit with the whole data. After understanding the “big picture,” we would like to find more models that could better fit different categories of process flow.

From the provided event log, we extracted unique paths and calculated the frequency with which is each path present in the data. Figure 4 depicts the most general process flow of the provided data. The primary process flows start with the purchasing document, creating the invoice by vendor, goods receipt, invoice receipt message and in the end, clearing the invoice. To measure how well the depicted process generalizes the actual process flow, we have counted all process instances and calculated the percentage of the cases that are going through each



Fig. 4. Generalized process flow.

path. As shown in the figure, even though the first and last event can be good generalized, with over 70% of instances going through them, the other process activities are describing only the smaller part of the process flow.

The process flows that are generated from the provided data have an unproportionally high amount of process variants. In total, there are 11 824 different flows the process can follow. That means that on average, every 21 process instances have different process flow. Consequently, a "spaghetti-like" process model results from the full, unfiltered event log. This unstructured process visualization is incomprehensible and cannot be used in operational process analysis.

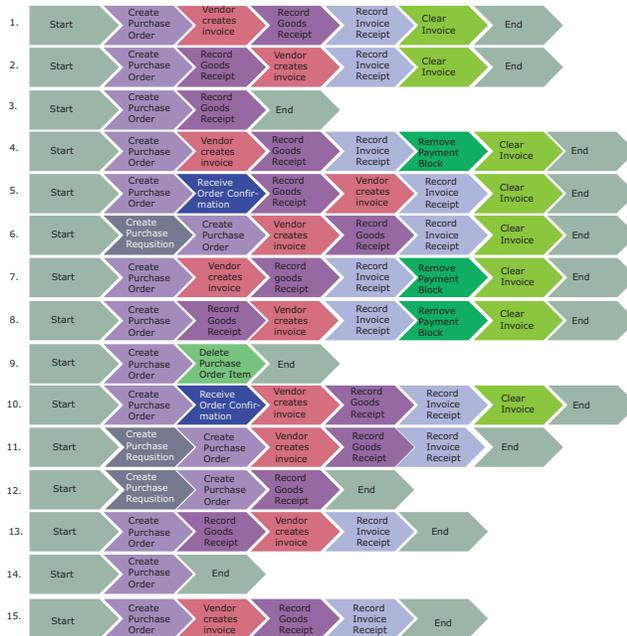
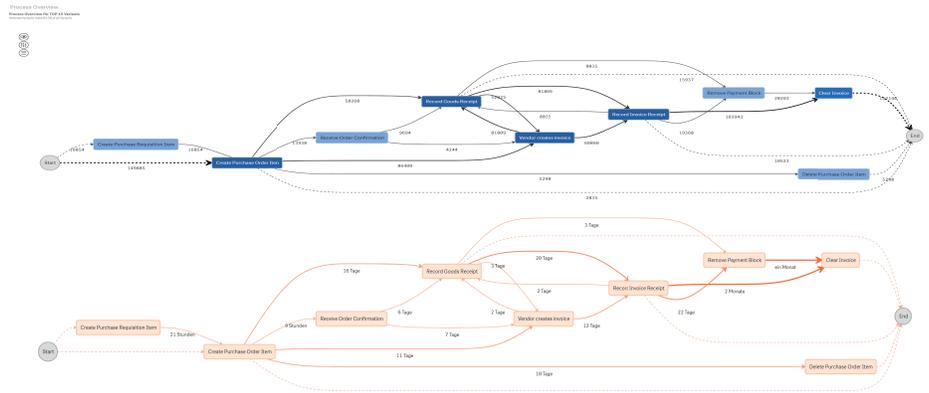


Fig. 5. Top 15 most frequent process variants.

To get a better understanding of the process, we have utilized the L\* approach for Spaghetti processes [3]. Consequently, we have identified 15 most frequent paths and calculated some of the performance metrics. Those most frequent

sequences are depicted in Fig. 5. The variants are grouped according to the count of instances that go through those flows. In the following chapters, we will use the numbers of paths for describing their properties. The generalized process flow shows that although the process map fits more than 66% of all process cases, it shows only about 50% of all events and less than 1 per cent of all process variants. Moreover, even if it shows the "birds-eye" view of the company, it does not help to answer questions under which conditions do the bottlenecks, or undesired paths take place.



**Fig. 6.** Process Flow with frequency and duration for top 15 most frequent process variants.

Figure 6 depicts the process models extracted from the variants described in table 2. The paths between most of the tasks take from 1 day to maximal two weeks. The biggest bottlenecks are between invoice receipt and removing of the payment block and invoice clearing activities. We can also see that there are 158. 606 Goods Receipt Messages, 142.669 Vendor Invoices but only 132. 146 cleared invoices. The analysis of the processes for singular months shows that whereas the general process model fitted from 70 - 80% of data in January to August 2018, the same model in September - December 2018 was applicable for only about 50% of all cases. In those months, almost 50% of all processes started with creating a purchase requisition item instead of purchase item. Likewise, the average cycle time of the process drops to 42 days instead of the 67 days.

Table 2 shows the top 15 process variants and some of the performance indicators for them. The top 15 Variants can describe ca. 65% of all process instances. So although the primary process flow variant fits every fifth purchase document; still, all the lower variants explain only the small percentage of the given data. The most straightforward process variants with the smallest amount of activities are also the ones with the shortest throughput.

**Table 2.** Process flow and performance indicators.

Pattern	# and % of Occurrences	Average (days)	St. dev(days)
1.	50.286 / 20%	70d 13h	33d 20h
2.	30.798 / 12%	84d 13h	33d 10h
3.	12.214 / 5%	22d 8h	21d 4h
4.	11.383 / 5%	91d 13h	41d 5h
5.	9.694 / 4%	93d 8h	20d 7h
6.	8.921 / 4%	60d 16h	23d 23h
7.	8.835 / 4%	77d 14h	39d 10h
8.	7.985 / 3%	102d 9h	34d 15h
9.	5.298 / 2%	9d 23h	25d 4h
10.	4.244 / 2%	76d 10h	23d 2h
11.	4.210 / 2%	38d 2h	19d 21h
12.	3.723 / 1%	26d 9h	20d
13.	3.548 / 1%	33d 8h	21d 23h
14.	2.835 / 1%	0d	0d
15.	2.765 / 1%	43d 3h	33d 12h

## 4 Challenge questions

Up until that point, we have described the general characteristics of the data and discovered the most widespread process model. In the next chapters, we want to answer the challenge questions and try to build the framework of performance indicators.

### 4.1 Q1: Collection of process models

Is there a collection of process models which together properly describe the process in this data. Based on the four categories above, at least 4 models are needed, but any collection of models that together explain the process well is appreciated. Preferably, the decision which model explains which purchase item best is based on properties of the item.

In the previous chapter, we have discovered and described model, which corresponds with ca. 65% of all cases. However, the generalized process model fits only the purchase documents that were created between January and August 2018. The standard deviations of throughput for different process variants were also relatively high, and the cycle times between them differed substantially.

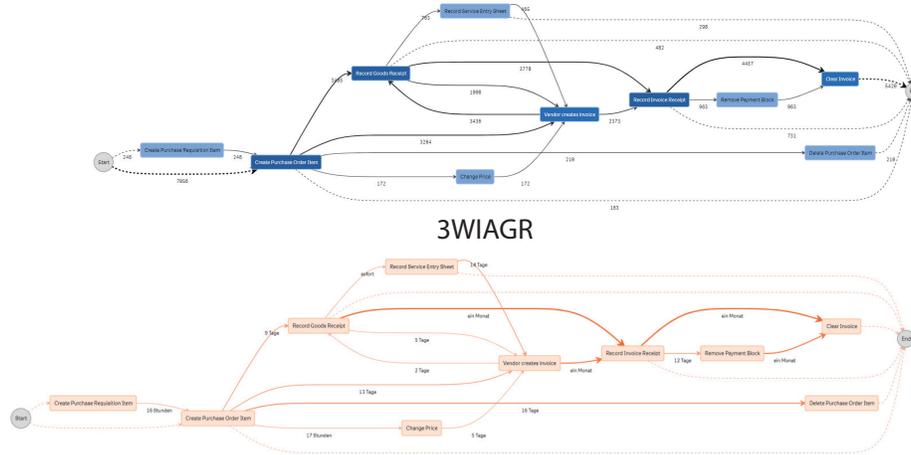
In the next step, we compared the general process model created in the previous chapter with all four process flows described in the challenge. As shown in table 3, the 3WIAGR and 3WIBGR document flow generally fit with the primary process model, 2W and especially CONS (Consignment process) does not many similarities with our previous process overview. Consequently, as in the previous chapter, we have separated the most common variants for all four case item types to discover the best fitting process models.

Figure 7 depicts the process model for 3WIAGR documents. The main difference between the general process and this picture are the additional tasks: record service entry sheet and change price. The throughput of the purchase

**Table 3.** Most frequent process variants for main flow types.

Pattern	% of Occurrences				
	All cases	3WIAGR	3WIBGR	2W	CONS
1.	50.286 / 20%	22%	15%	0%	0%
2.	30.798 / 12%	13%	9%	0%	0%
3.	12.214 / 5%	1%	3%	61%	0%
4.	11.383 / 5%	5%	3%	0%	0%
5.	9.694 / 4%	4%	0%	0%	0%
6.	8.921 / 4%	4%	2%	0%	0%
7.	8.835 / 4%	4%	0%	0%	0%
8.	7.985 / 3%	3%	2%	0%	0%
9.	5.298 / 2%	2%	1%	2%	0%
10.	4.244 / 2%	2%	0%	0%	0%
11.	4.210 / 2%	2%	0%	0%	0%
12.	3.723 / 1%	1%	1%	12%	0%
13.	3.548 / 1%	2%	2%	0%	0%
14.	2.835 / 1%	1%	1%	2%	0%
15.	2.765 / 1%	1%	1%	0%	0%

orders of 3WIAGR type is 70 days, and each process has, on average, 21 activities. 42% of all SRM Orders are included in that category. The documents are created only by companyID\_0000 and companyID\_0001. About 10% of all invoices of that type are not cleared.



**Fig. 7.** 3-way match, invoice after GR.

Figure 8 depicts the process model with 15 most common variants for 3WIBGR documents. About 88% of all purchase orders belong to that category. T. The throughput of the purchase orders of 3WIAGR type is 73 days. In opposite to 3WIAGR documents the typical process in this category has only 5 activities.



3WIAGR and 3WIBGR category. The typical process has 5 activities. There are no goods receipt messages in that process type and no SRM tasks. The biggest group of case spend area are Real estate and others purchases. There are no product-related documents in this process flow.

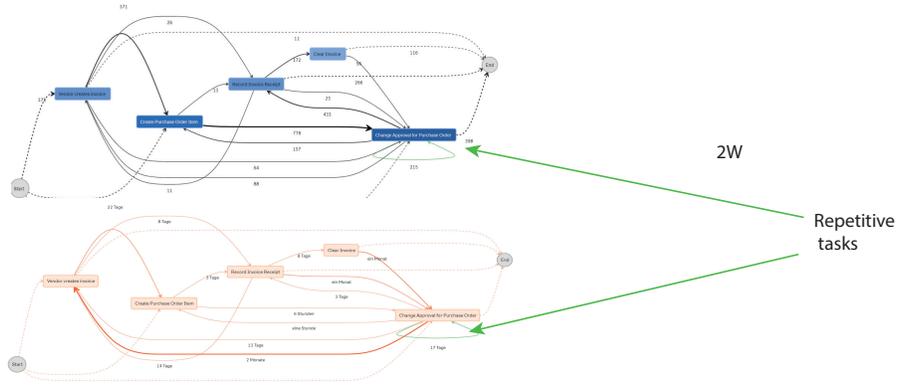


Fig. 10. 2 way match.

With the general model, we could answer the question of what happens in the process but not why and under what circumstances does it happen.

4.2 Q2: Invoicing process throughput

What is the throughput of the invoicing process, i.e. the time between goods receipt, invoice receipt and payment (clear invoice)? To answer this, a technique is sought to match these events within a line item, i.e. if there are multiple goods receipt messages and multiple invoices within a line item, how are they related and which belong together?.

In our general analysis, we used the compliance flag in to compare the first amount of purchase document (value at the creation) with the sum of values connected to goods receipt items and the same for the invoice receipt messages. That value matches for almost 85% of all cases and 61 of the total purchase volume. Similarly, for the second question, we have loaded only the events from the invoicing process. From the filtered cases, only about 20% (51.497) of purchasing documents do not have matching values. Figure 11 shows the attributes for non-matching and matching values. Figure 11 depicts the typical attributes of the purchasing documents. Both categories have mostly 3WIBGR type documents. The most significant difference is by not matching documents the proportion of product-related purchases, which is much bigger than by matching documents.

Figure 12 shows the comparison of performance metrics for both matching- and not matching the category of purchase documents. The documents that do

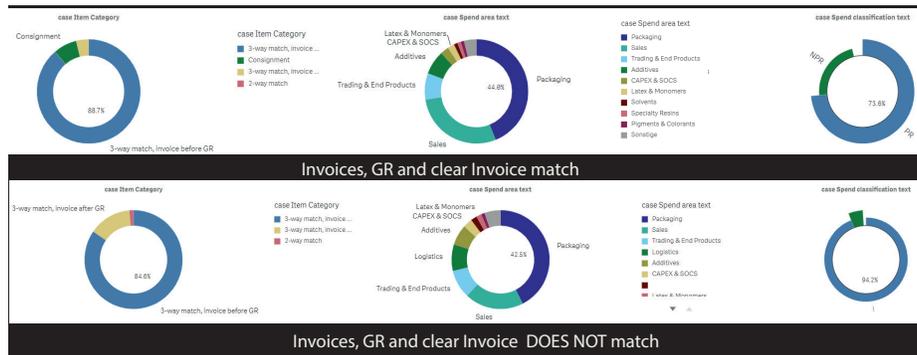


Fig. 11. Comparison of matching and not matching Purchasing Documents.

not match have usually more activities pro process (more repetitive tasks) but at the same time, the cycle time is much shorter. The amount for each purchase documents is for that category bigger than for those with matching values.

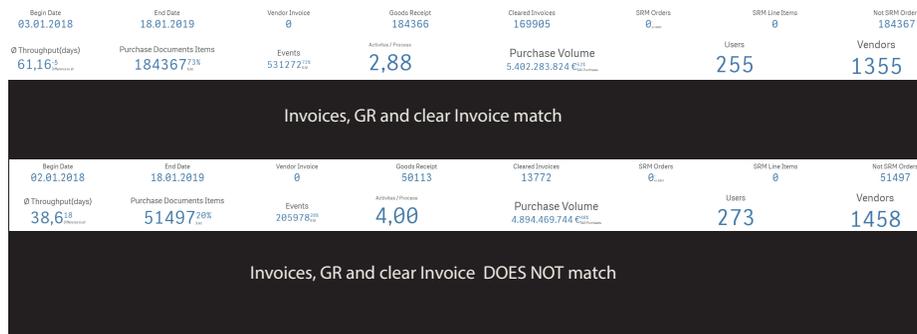


Fig. 12. Comparison of matching and not matching Purchasing Documents.

Figure 13 shows the top 10 vendors for which the values do not match. All the documents for those vendors belong to the non-product related category. Most of the purchases are involving logistics or marketing. Although there are only 843 purchase document items, they add up to 8% of the total purchase volume. All of those documents belong to the 3WIAGR category.

Figure 14 compares process models for both matching and not matching categories. The process models show much more repetitions for not matching documents. The repetitive tasks are on average longer than in those matching cases. To compare the two categories, we have classified all documents in cat-

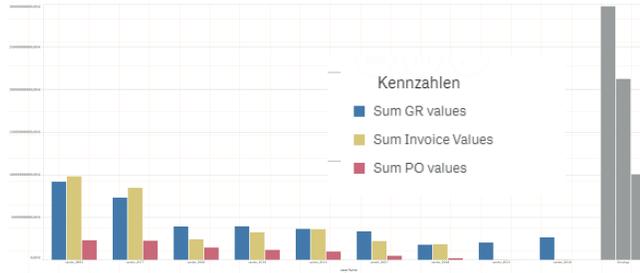


Fig. 13. Top 10 Vendors for not matching documents.

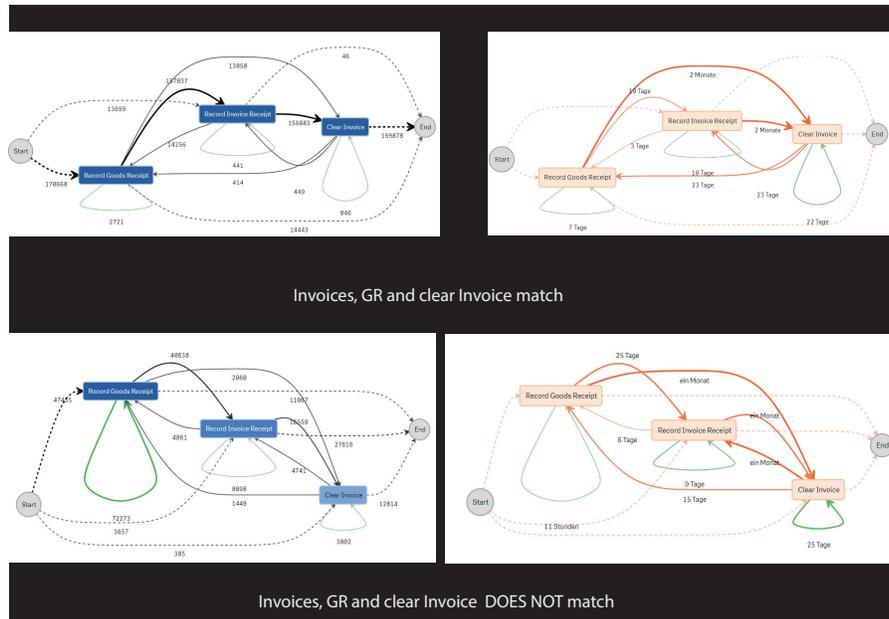
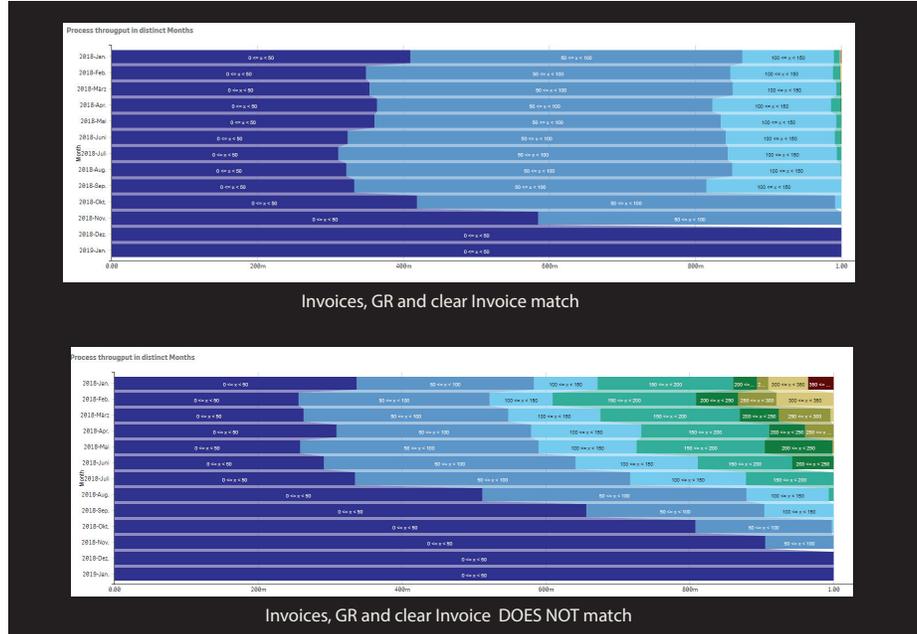


Fig. 14. Comparison of process models for matching and not matching Purchasing Documents.

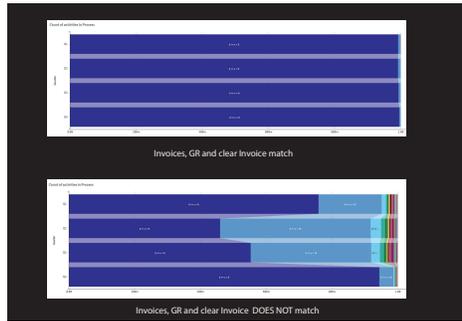
egories comparing their throughput time. Figure 15 shows the comparison of those two types of documents based on throughput clusters. The most frequent cases have cycle times from 0 to 50 days (dark blue). Especially at the beginning of 2018, there is a clear difference in throughput times between those documents.



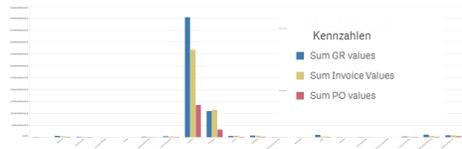
**Fig. 15.** Comparison of cycle times for matching and not matching Purchasing Documents.

Similarly to throughput times, we have divided all documents according to the number of events for each process. Figure 16 shows a comparison of those groups. The cases in which there are matching values have at most 5 activities per process, and those with outstanding amounts have more activities and repetitive tasks.

Figure 17 shows the spend areas for which values in the invoice, goods receipt and payment does not match. Especially in the logistics process, the amount of goods receipt is much higher than invoice and payment values. Other typical categories for which the singular documents do not match are marketing, workforce services and packaging.



**Fig. 16.** Case spend areas for not matching Purchasing Documents.



**Fig. 17.** Comparison of matching and not matching Purchasing Documents.

### 4.3 Q3: Outliers and deviations

Finally, which Purchase Documents stand out from the log? Where are deviations from the processes discovered in (1) and how severe are these deviations? Deviations may be according to the prescribed high-level process flow, but also with respect to the values of the invoices. Which customers produce a lot of rework as invoices are wrong, etc.?

As a part of the data load, we have created different flags for the cases that are, in some way deviating from the normal process flow or could cause in some way compliance or rules violation. We have found that about 34% of all cases have different deviations or suspect behaviour. The most frequent process deviation is rapid goods delivery. There are about 6 thousand cases where the receipt of the goods follows other tasks only after a few minutes. This kind of behaviour is typical for vendors 0458, 0522,0228 and 0525. The affected spend areas are logistics, capex & socs and sales.

The other suspected behaviour are the cases where purchase order item is created three or more weeks after the vendor invoice. This late Purchase Order creation is typical for about 3 thousand documents. This behaviour is found mainly at Trading & End products. In about 2400 cases the price of the purchase order is changed after an invoice has been issued. It is usually the case for packaging, sales and additives purchases.

## 5 Conclusion

As a part of BPI Challenge 2019, we have analyzed the purchasing process for a large multinational company. The decentralized process was not only complicated but did involve many different subprocesses. On the one hand, the big part of purchasing involved non-product related purchases. Especially the processes involving marketing or logistics contained many deviations and untypical behaviour.

## References

1. van Dongen, B.F., Dataset BPI Challenge 2019. 4TU.Centre for Research Data. <https://doi.org/10.4121/uuid:d06aff4b-79f0-45e6-8ec8-e19730c248f1>
2. Telgen, J. and de Boer, L. 1995. Developments in purchasing of non-production items. Proceedings 4th IPSERA Conference, Eindhoven, 1-8.
3. W. M. P. van der Aalst, "Process mining: discovering and improving Spaghetti and Lasagna processes," 2011 IEEE Symposium on Computational Intelligence and Data Mining (CIDM), Paris, 2011, pp. 1-7. <https://doi.org/10.1109/CIDM.2011.6129461>