

# Process Mining With Jupyter Notebook And Disco





**Python** have the best quality of documentation.

The maturity and stability of the fundamental numerical libraries are well known.

**pandas** is a Python library of data structures and statistical tools.

**Jupyter** is the GUI tool set that widely use for the Python open source community for Data pre-processing.



## The Data Set - **BPI Challenge 2019**

The data was recorded for the execution of **Purchase-To-Pay** processes.

- It was recorded possibly from an ERP system.
- The log contains 1,595,923 events, with **22 columns**.
- Character set - Windows-1252 or CP-1252

1595923 rows × 22 columns

```
df = pd.read_csv("BPI2019.csv", encoding='cp1252')
```

# Discover : Case ID

#	Column	Non-Null Count	Dtype
0	eventID	1595923 non-null	int64

```
In [6]: df['eventID'].unique()
Out[6]: array([ 65781719105536, 65777424138241, 65777424138240, ...,
1009441868611588, 273679611068428, 77635828842576], dtype=int64)

In [7]: pd.value_counts(df['eventID'])
Out[7]: 697502688870404    1
1058967136501761    1
610233248382980    1
211767657496581    1
77249281785859    1
..
214318868070410    1
796652008898560    1
911705592823808    1
682143885819904    1
0    1
Name: eventID, Length: 1595923, dtype: int64

In [8]: count = df['eventID'].value_counts()

In [9]: (count == 1).value_counts()
Out[9]: True    1595923
Name: eventID, dtype: int64
```

**eventID** can not be a **case ID** for Process Discovery

```
In [10]: df.loc[df['eventID'] == 697502688870404, :]
Out[10]:
```

	eventID	case Spend area text	case Company	Document Type	case Sub spend area text	case Purchasing Document	case Purch. Doc. Category name	case Vendor	case Item Type	case Item Category	...	case Name	case GR-Based Inv. Verif.	case
1505633	697502688870404	Packaging	companyID_0000	Standard PO	Labels	4508048579	Purchase order	vendorID_0120	Standard	3-way match, invoice before GR	...	vendor_0119	False	

1 rows x 22 columns

## Discover : Case ID

#	Column	Non-Null Count	Dtype
0	eventID	1595923 non-null	int64
1	case Spend area text	1579629 non-null	object
2	case Company	1595923 non-null	object
3	case Document Type	1595923 non-null	object
4	case Sub spend area text	1579629 non-null	object
5	case Purchasing Document	1595923 non-null	int64
6	case Purch. Doc. Category name	1595923 non-null	object
7	case Vendor	1595923 non-null	object
8	case Item Type	1595923 non-null	object

### case Purchasing Document

- No null value
- Int64
- But no flow **Activity**

**case concept:name** Is the perfect match

- Same PO
- Good flow of Activity
- Same Amount (EUR)
- Continuous timestamp

```
In [19]: case_id = df['case concept:name'].replace('_', '', regex=True)
```

case concept:name	case Goods Receipt	event User	event org:resource	event concept:name	event Cumulative net worth (EUR)	event time:timestamp
4508048579_00130	True	user_057	user_057	Create Purchase Order Item	418.0	10-08-2018 16:20:00.000
4508048579_00130	True	user_034	user_034	Record Goods Receipt	418.0	29-08-2018 11:17:00.000
4508048579_00130	True	NONE	NONE	Vendor creates invoice	418.0	30-08-2018 23:59:00.000
4508048579_00130	True	user_012	user_012	Record Invoice Receipt	418.0	31-08-2018 12:26:00.000
4508048579_00130	True	user_023	user_023	Remove Payment Block	418.0	13-12-2018 08:03:00.000
4508048579_00130	True	user_002	user_002	Clear Invoice	418.0	13-12-2018 14:08:00.000

```
In [23]: df['case concept:name'] = pd.to_numeric(df['case concept:name'])
```

event time:timestamp

```
In [14]: pd.to_datetime(df['event time:timestamp']).head()
```

```
Out[14]: 0    1948-01-26 23:59:00
         1    1948-01-26 23:59:00
         2    1948-01-26 23:59:00
         3    1948-01-26 23:59:00
         4    1948-01-26 23:59:00
         Name: event time:timestamp, dtype: datetime64[ns]
```

```
In [27]: df.loc[df['timestamp'] < '2017-12-30 00:00:00', 'timestamp']
```

```
Out[27]: 0    1948-01-26 23:59:00
         1    1948-01-26 23:59:00
         2    1948-01-26 23:59:00
         3    1948-01-26 23:59:00
         4    1948-01-26 23:59:00
         ...
        290  2017-12-29 23:59:00
        291  2017-12-29 23:59:00
        292  2017-12-29 23:59:00
        293  2017-12-29 23:59:00
        294  2017-12-29 23:59:00
         Name: timestamp, Length: 295, dtype: datetime64[ns]
```

```
In [25]: df['timestamp'].dt.year.value_counts()
```

```
Out[25]: 2018    1550468
         2019     45135
         2017      223
         2008       45
         2001       22
         1948       10
         1993        9
         2016        6
         2015        3
         2020        2
         Name: timestamp, dtype: int64
```

It was a Object type data.

For the search and manipulation we change it to **datetime64**

- Purchase orders submitted mostly in 2018.
- Some purchase orders start from 1948 (only 295 out of 1,595,923)
- These 295 purchase orders are outlier
- We delete those process



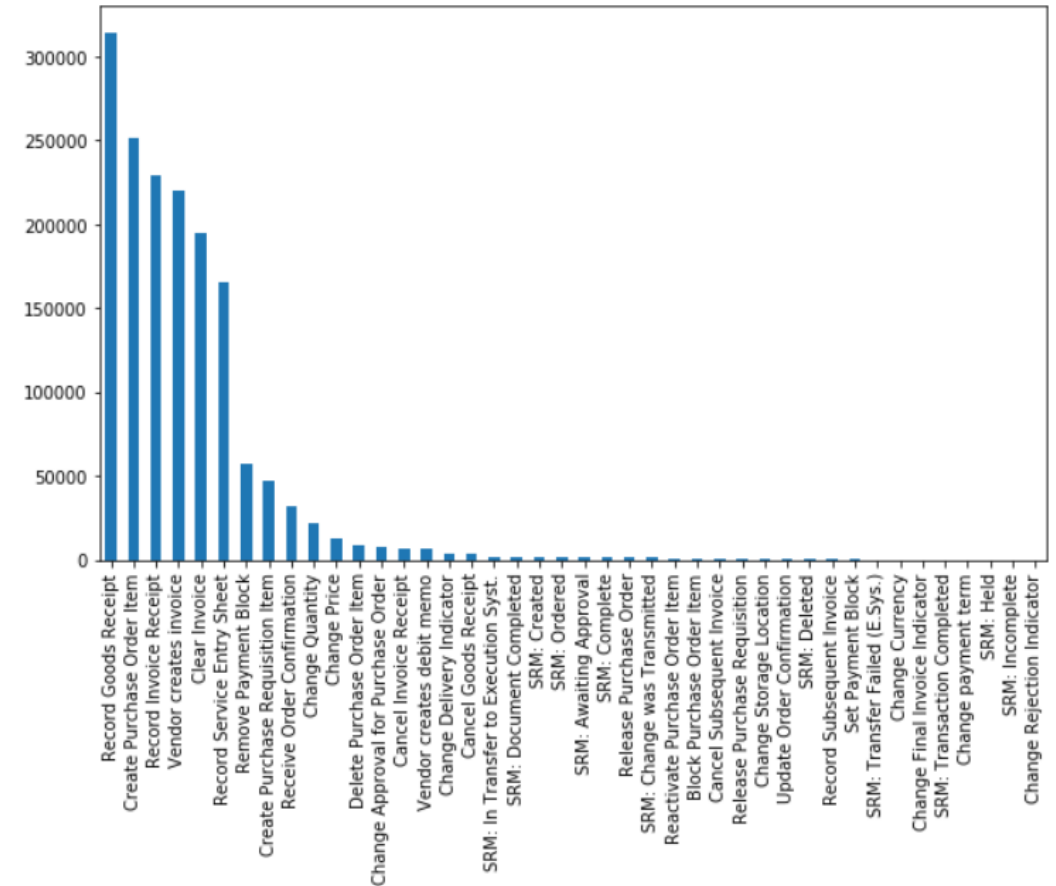
## Discover : Activity

event concept:name

Activity	Frequency
Record Goods Receipt	314097
Create Purchase Order Item	251734
Record Invoice Receipt	228760
Vendor creates invoice	219919
Clear Invoice	194393
Record Service Entry Sheet	164975
Remove Payment Block	57136
Create Purchase Requisition Item	46592
Receive Order Confirmation	32065
Change Quantity	21449
Change Price	12423
Delete Purchase Order Item	8875
Change Approval for Purchase Order	7541
Cancel Invoice Receipt	7096
Vendor creates debit memo	6255
Change Delivery Indicator	3289
Cancel Goods Receipt	3096

```
In [7]: df['event concept:name'].value_counts().plot.bar(figsize=(10,6))
```

```
Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x232a0c56588>
```



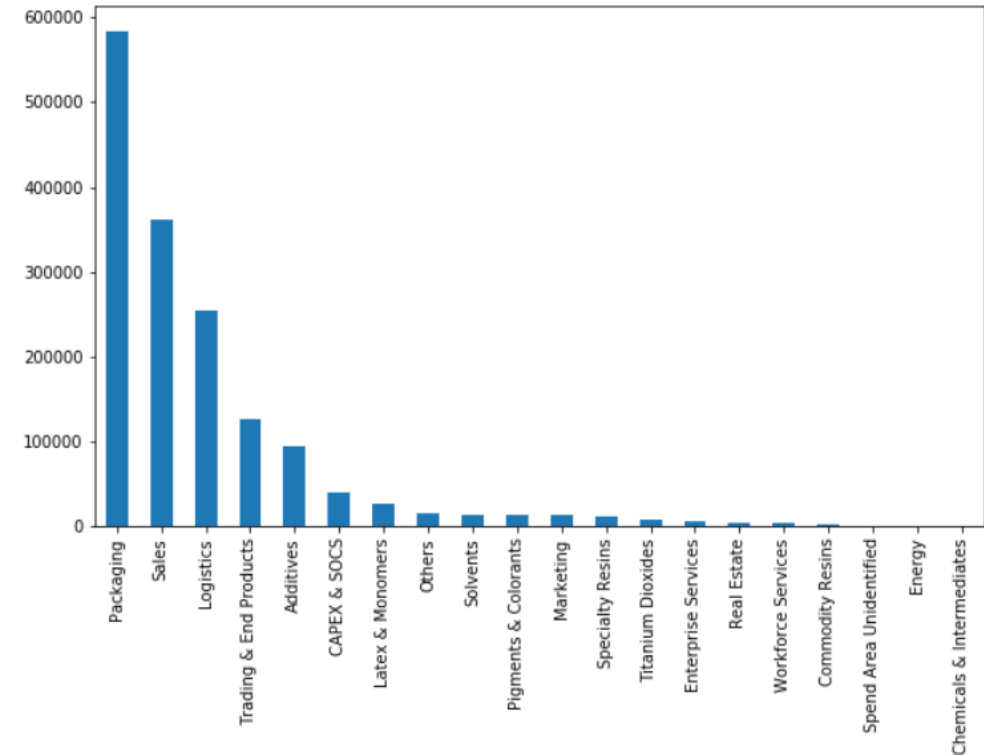
## Discover : Activity

## case Spend area text

Activity	Frequency
Packaging	583981
Sales	360774
Logistics	253565
Trading & End Products	126756
Additives	95499
CAPEX & SOCS	40074
Latex & Monomers	27007
Others	15419
Solvents	13889
Pigments & Colorants	13811
Marketing	12994
Specialty Resins	12469
Titanium Dioxides	7453
Enterprise Services	5957
Real Estate	3824
Workforce Services	3441
Commodity Resins	2374

```
In [14]: df['case Spend area text'].value_counts().plot.bar(figsize=(10,6))
```

```
Out[14]: <matplotlib.axes._subplots.AxesSubplot at 0x232a1cc4808>
```





## Discover : User

```
In [7]: print('===case Vendor===')
print(df['case Vendor'].unique())
print('')

print('===case Name===')
print(df['case Name'].unique())
print('')

print('===event User ===')
print(df['event User'].unique())
print('')

print('===event org:resource ===')
print(df['event org:resource'].unique())
print('')
```

Four Columns  
case Vendor  
case Name  
event User  
event org:resource

```
===case Vendor===
['vendorID_0670' 'vendorID_0427' 'vendorID_0307' ... 'vendorID_1968'
 'vendorID_1973' 'vendorID_1974']

===case Name===
['vendor_0645' 'vendor_0415' 'vendor_0298' ... 'vendor_1892' 'vendor_1897'
 'vendor_1898']

===event User ===
['NONE' 'user_329' 'user_236' 'user_124' 'batch_03' 'batch_08' 'batch_04'
 'user_033' 'user_036' 'user_038' 'batch_00' 'user_043' 'user_045'
 'user_051' 'user_029' 'user_052' 'user_054' 'user_057' 'user_059'
 'user_060' 'user_064' 'user_066' 'user_068' 'user_070' 'user_072'
 'user_075' 'user_079' 'user_084' 'user_085' 'user_081' 'user_089'
 'user_091' 'user_092' 'user_095' 'user_097' 'user_100' 'user_103'
 'user_105' 'user_108' 'user_000' 'user_110' 'user_116' 'batch_06'
 'batch_02' 'user_154' 'user_086' 'user_113' 'user_118' 'user_122']
```

## Data quality issues

Data Quality Issues	BPI Challenge 2019
Incorrect Timestamps	X
Missing values	
Missing Events	X
Duplicate Tasks	X
Overlapping Activity Executions	X
Case Heterogeneity	X
Voluminous Data	X
Noisy Data-Outliers	X

## Final Data set

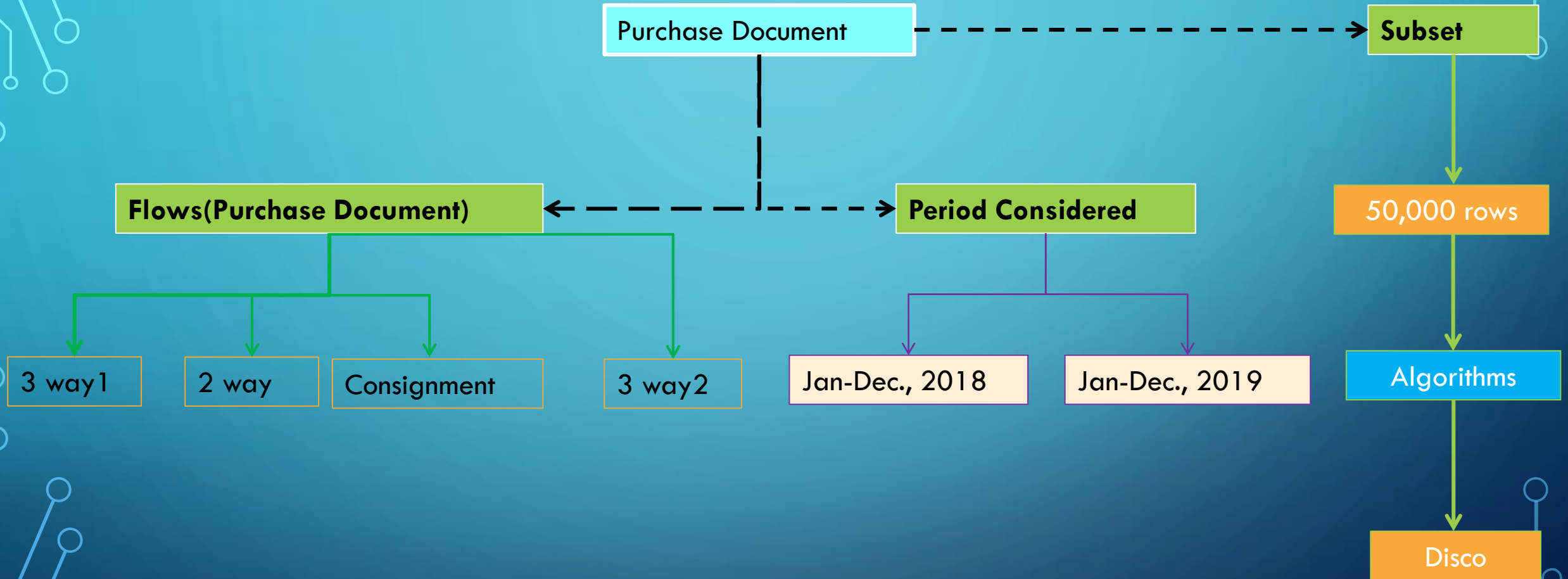
```
In [32]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1588570 entries, 0 to 1588569
Data columns (total 13 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   case Spend area text                 1572311 non-null object  
 1   case Sub spend area text             1572311 non-null object  
 2   case Vendor                          1588570 non-null object  
 3   case Item Type                       1588570 non-null object  
 4   case Item Category                   1588570 non-null object  
 5   case Spend classification text        1572311 non-null object  
 6   case Name                            1588570 non-null object  
 7   case GR-Based Inv. Verif.            1588570 non-null bool     
 8   case concept:name                    1588570 non-null int64  
 9   event User                           1588570 non-null object  
10   event concept:name                   1588570 non-null object  
11   event Cumulative net worth (EUR)      1588570 non-null float64 
12   timestamp                            1588570 non-null object  
dtypes: bool(1), float64(1), int64(1), object(10)
memory usage: 147.0+ MB
```

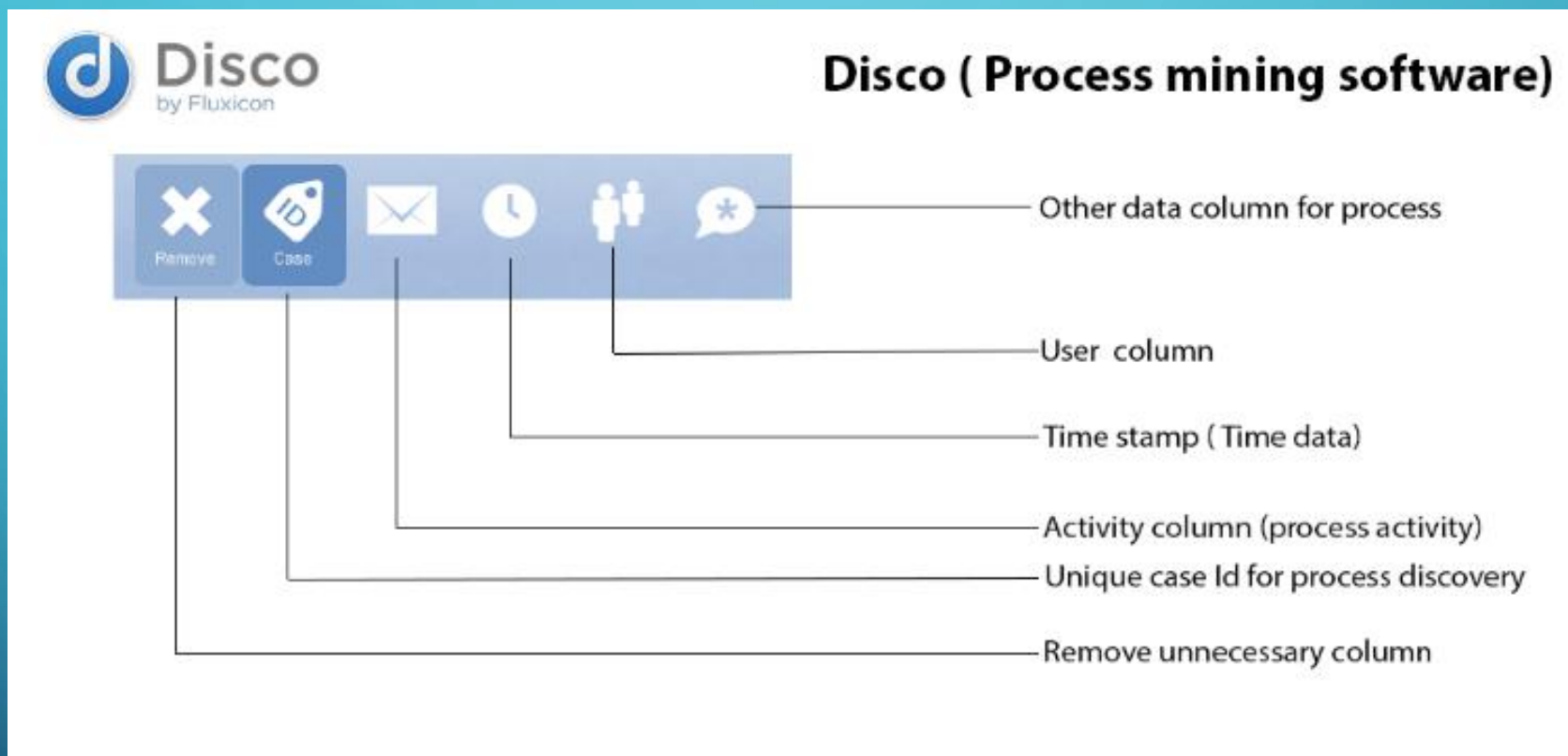
```
In [27]: df.to_csv('B03_done.csv', encoding='utf-8', index=False)
```

# Dataset

Focusing:



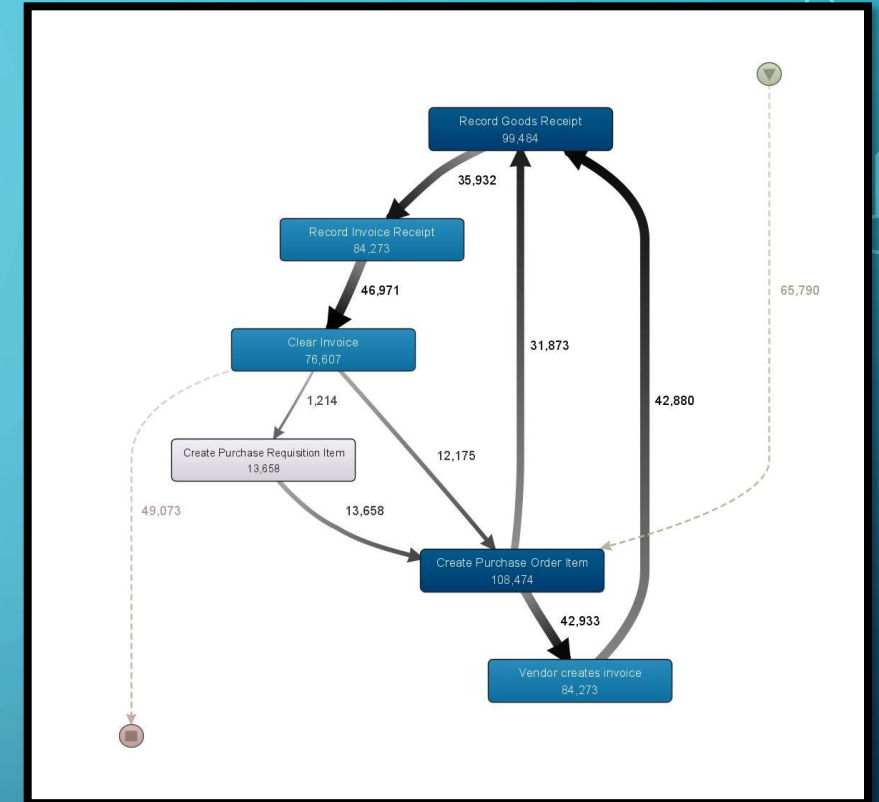
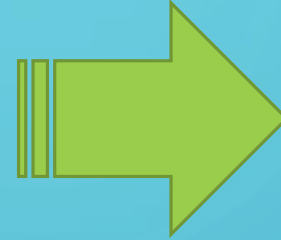
# DISCO Implementation



# The Data Set - BPI Challenge 2019

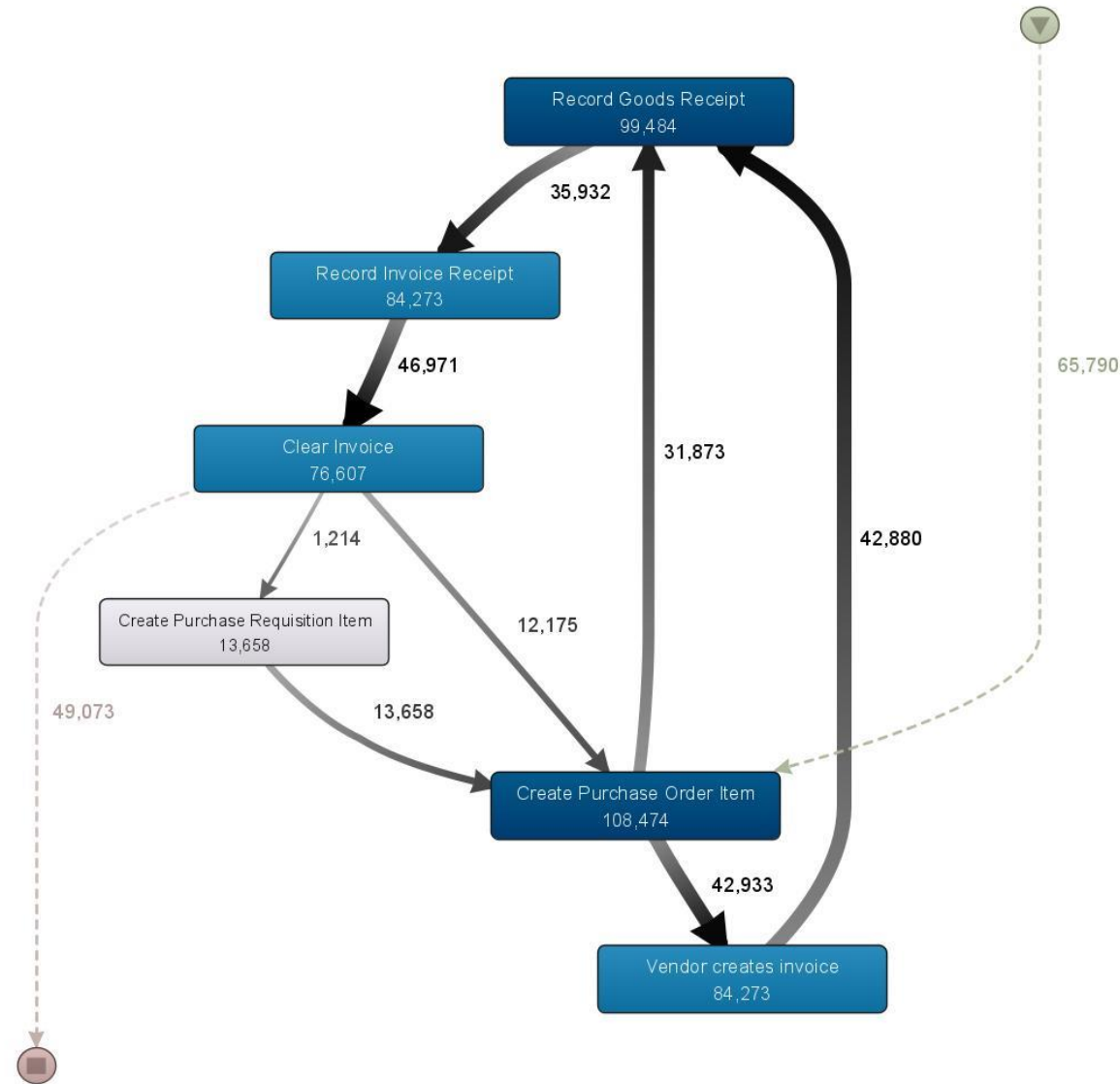
	eventID	case Spend area text	case Company	case Document Type	case Sub spend area text	case Purchasing Document	case Purch. Doc. Category name	case Vendor	case Item Type	case Item Category	...	case Name	case GR- Based Inv. Verif.	case Item
0	65781719105536	Sales	companyID_0000	Standard PO	Products for Resale	4507004931	Purchase order	vendorID_0670	Standard	3-way match, invoice before GR	...	vendor_0645	False	20 4507
1	65777424138241	Sales	companyID_0000	Standard PO	Products for Resale	4507004931	Purchase order	vendorID_0670	Standard	3-way match, invoice before GR	...	vendor_0645	False	10 4507
2	65777424138240	Sales	companyID_0000	Standard PO	Products for Resale	4507004931	Purchase order	vendorID_0670	Standard	3-way match, invoice before GR	...	vendor_0645	False	10 4507
3	65794604007424	Sales	companyID_0000	Standard PO	Products for Resale	4507004931	Purchase order	vendorID_0670	Standard	3-way match, invoice before GR	...	vendor_0645	False	50 4507
4	65794604007425	Sales	companyID_0000	Standard PO	Products for Resale	4507004931	Purchase order	vendorID_0670	Standard	3-way match, invoice before GR	...	vendor_0645	False	50 4507

## Pre-processing



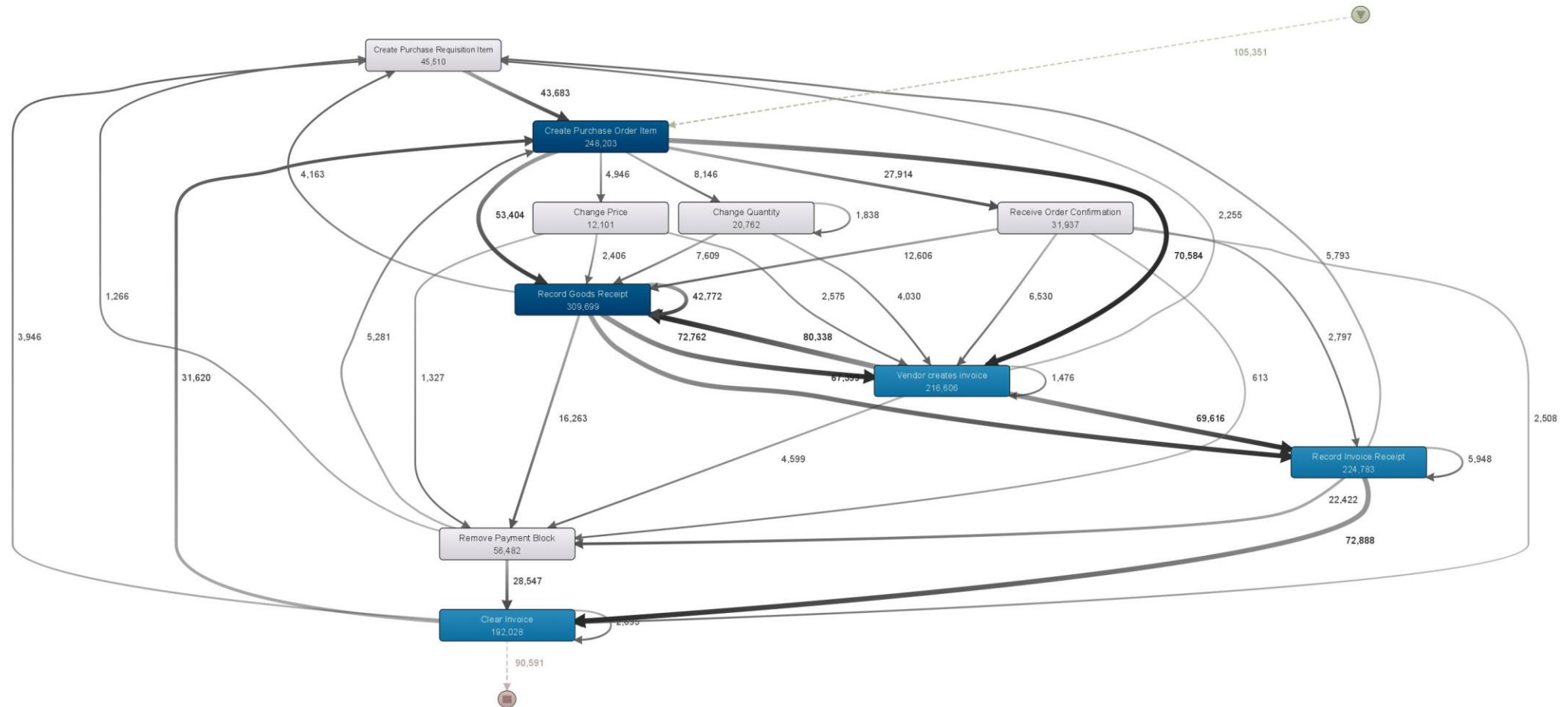
- Case ID
- Timestamp
- Activity
- Other fields.....

# Activity

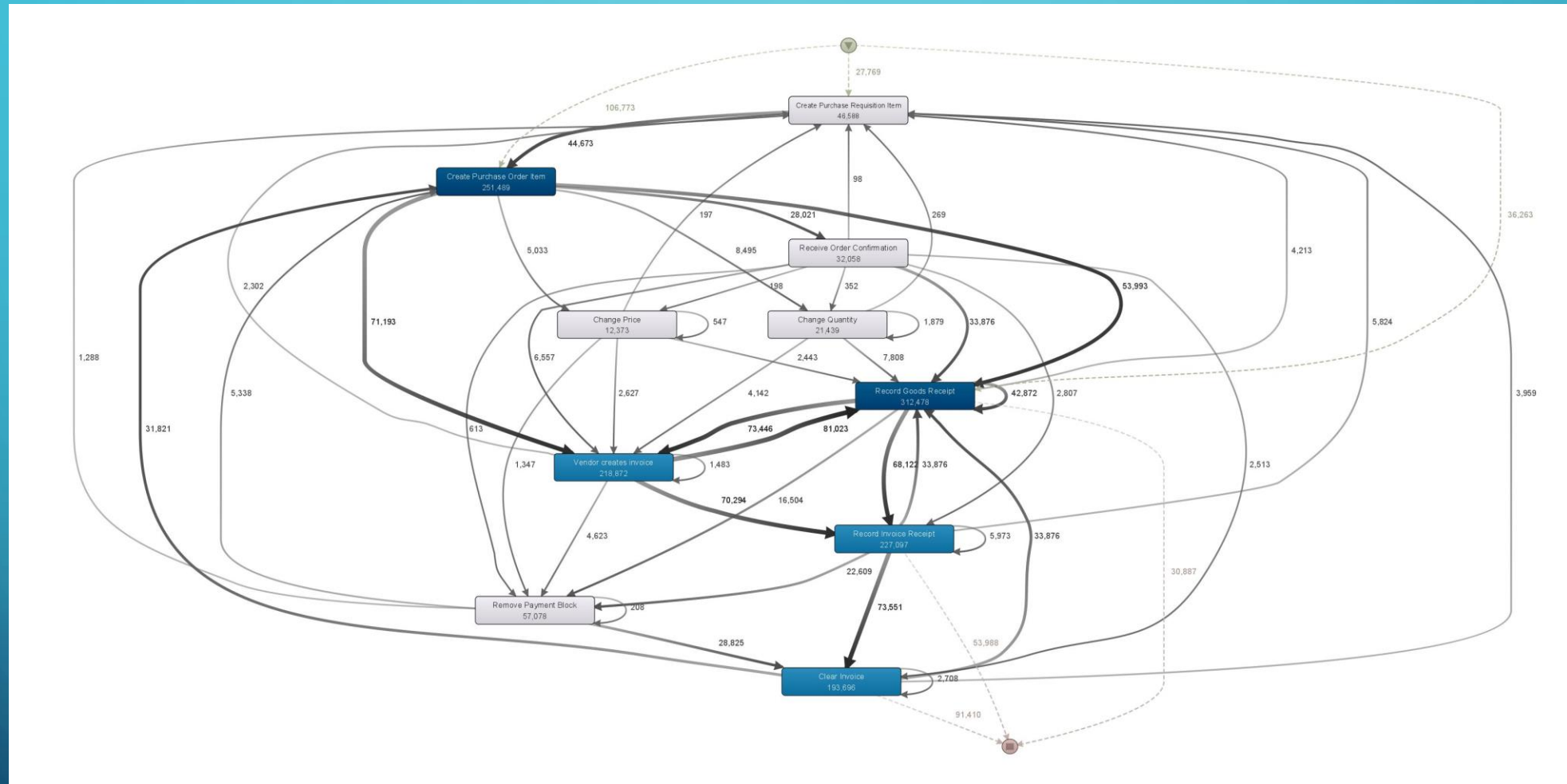




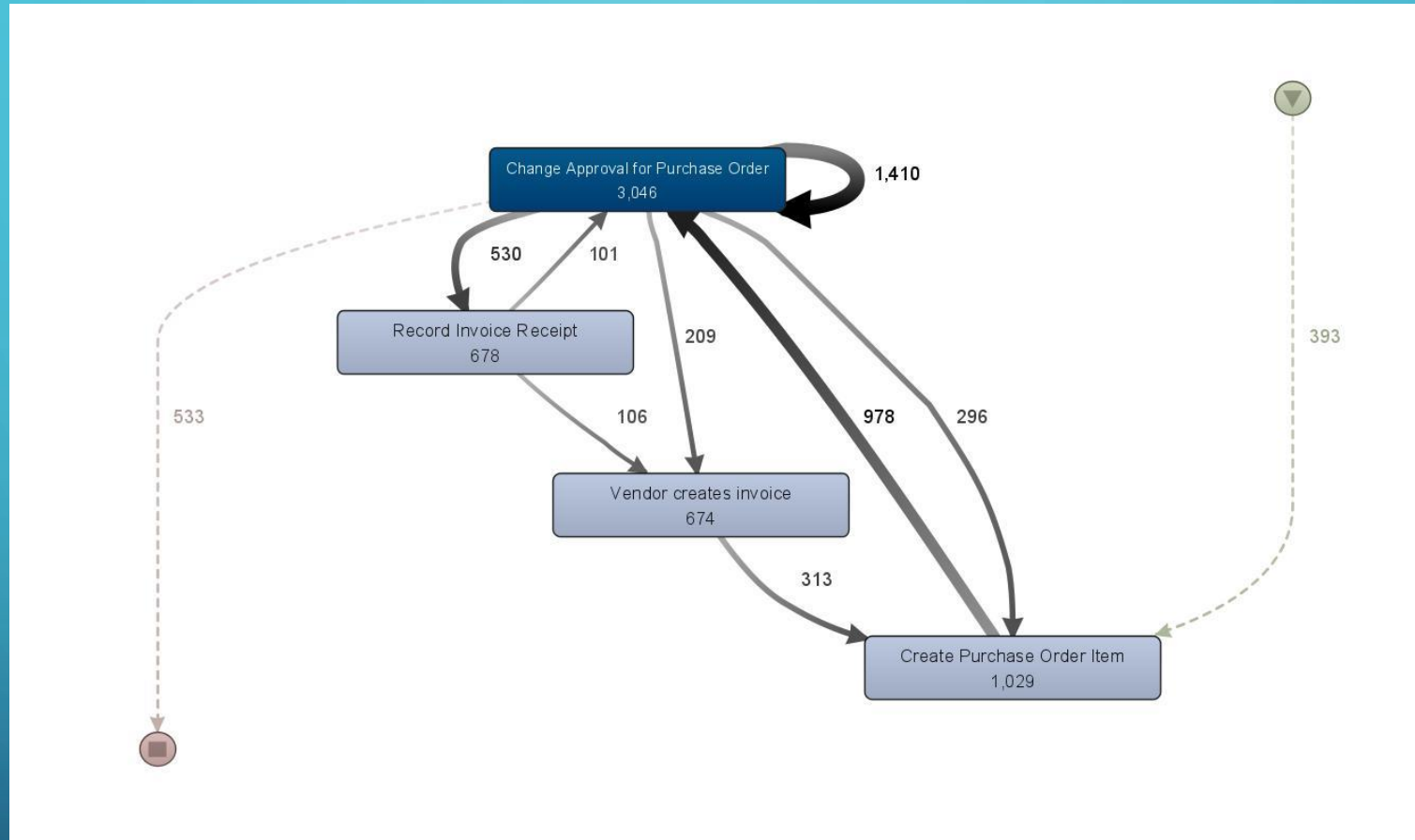
## Case Spend area text



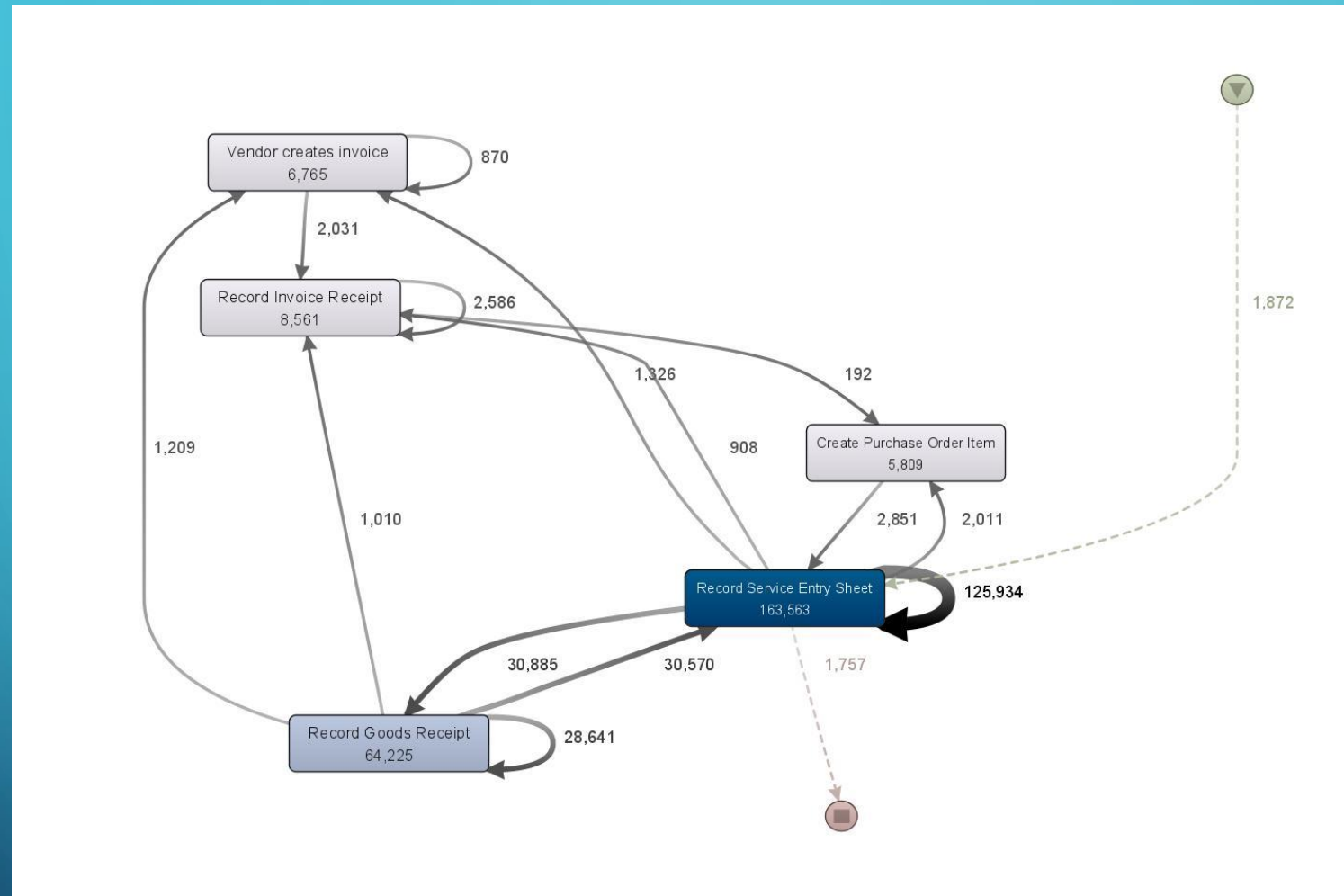
## Case Sub spend area text



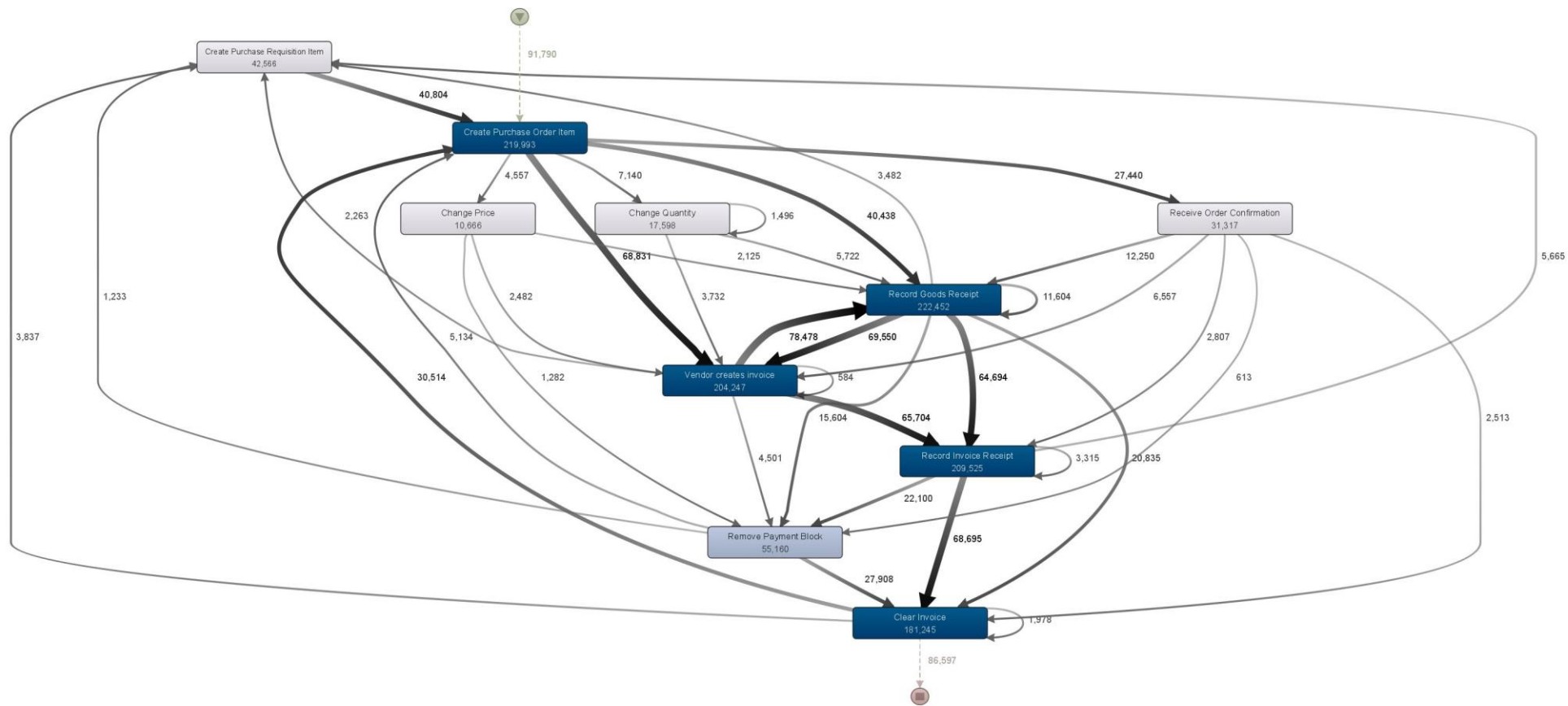
## Case Item Type - limit



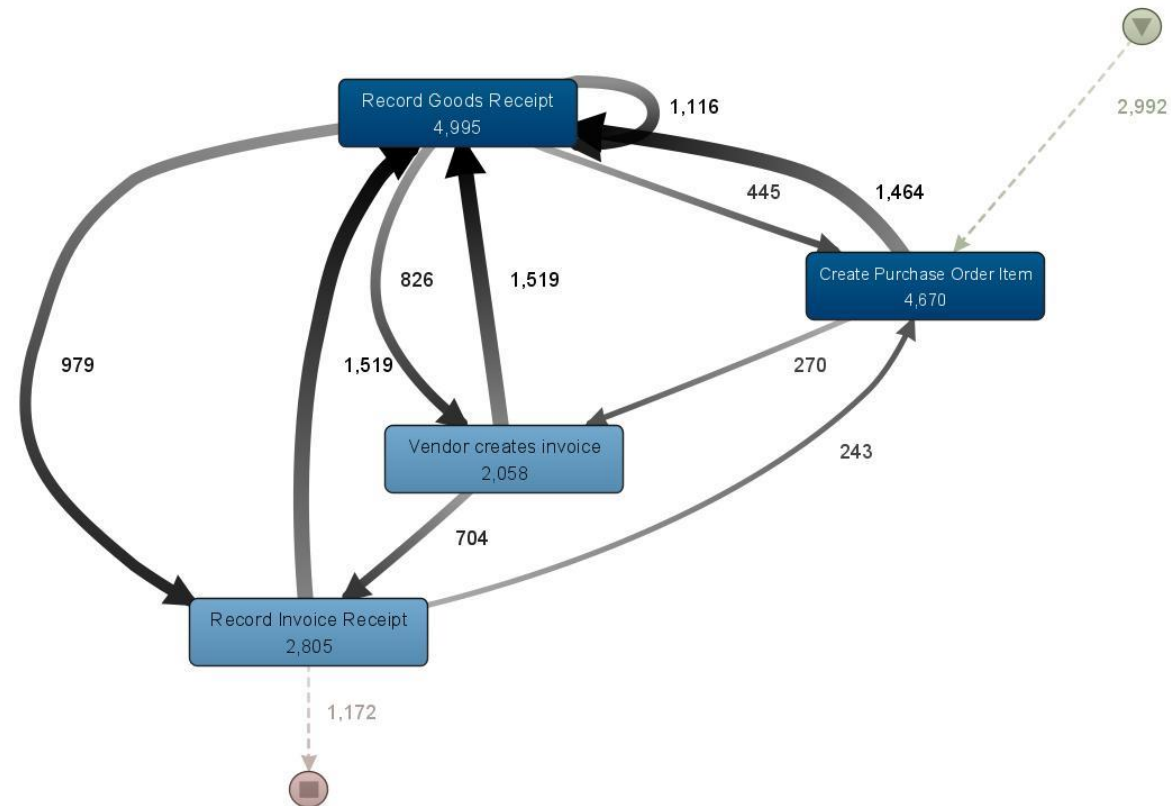
## Case Item Type - service



## Case Item Type - standard

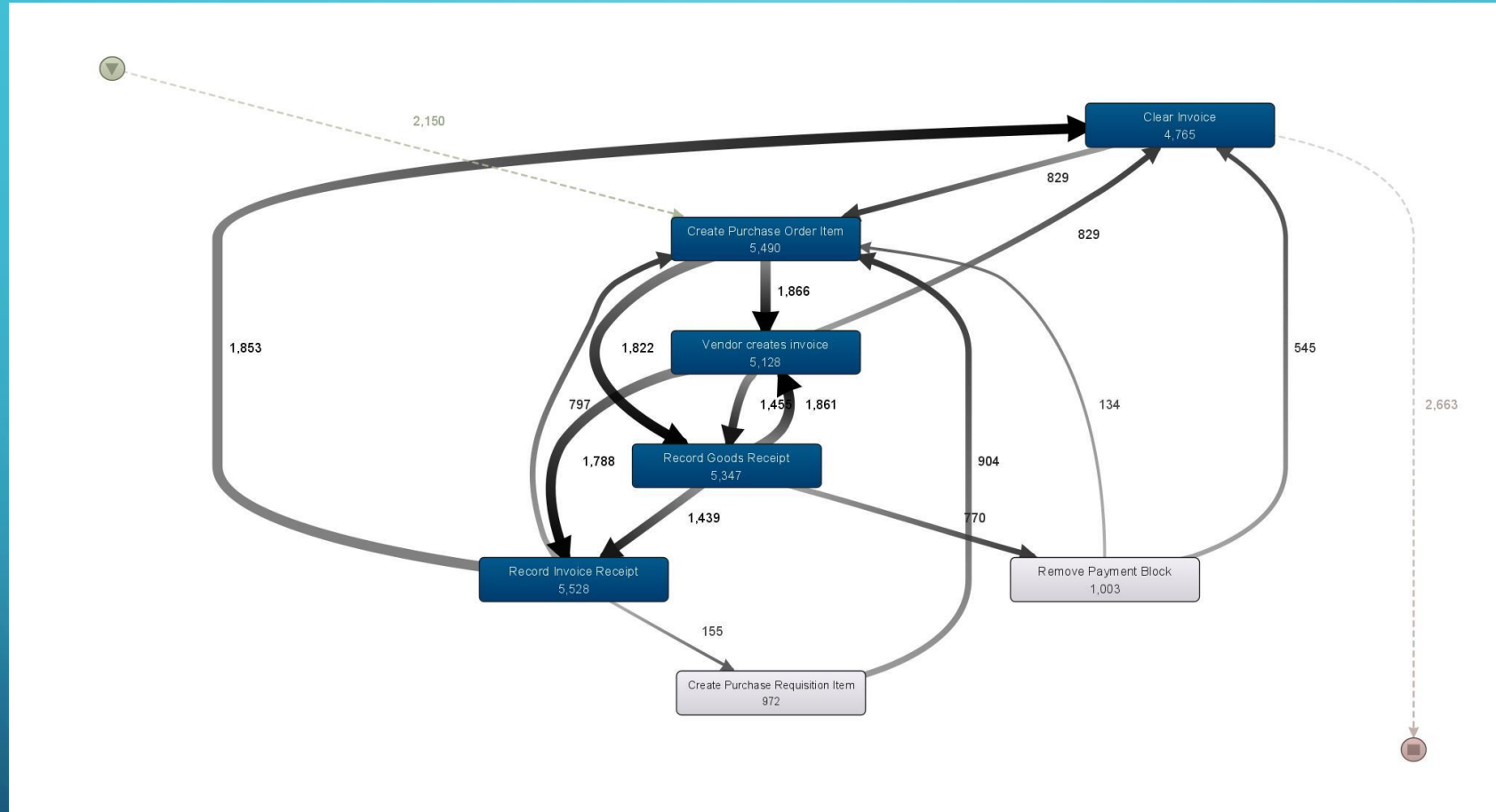


## Case Item Type - subcontract





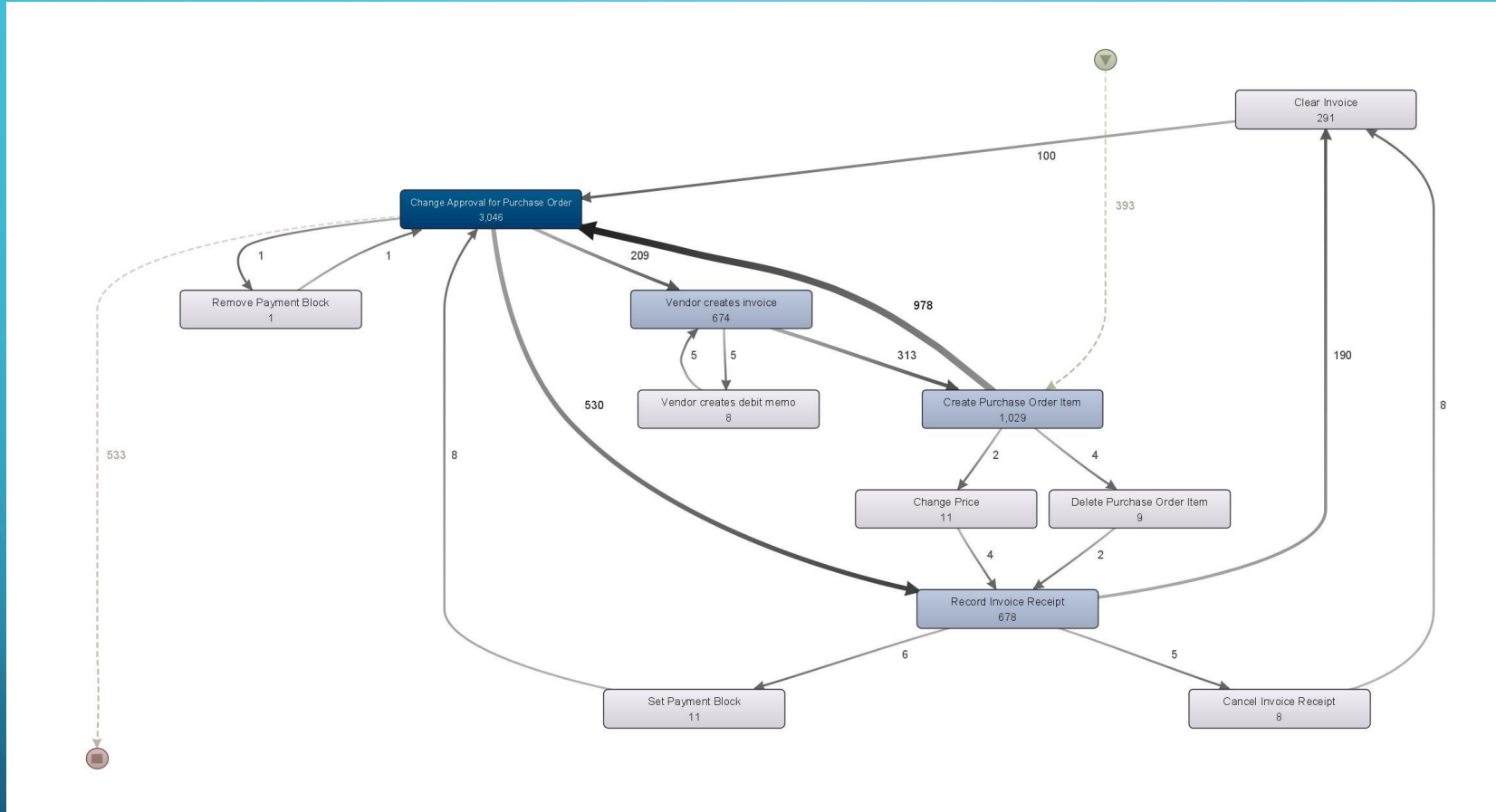
## Case Item Type - Third party



## Case item category – 2-way match

**(Values)** Invoice ----matches creation (PO)

**(Implication)** No separate goods receipt message required.  
GR-based flag and the Goods Receipt flags set to false



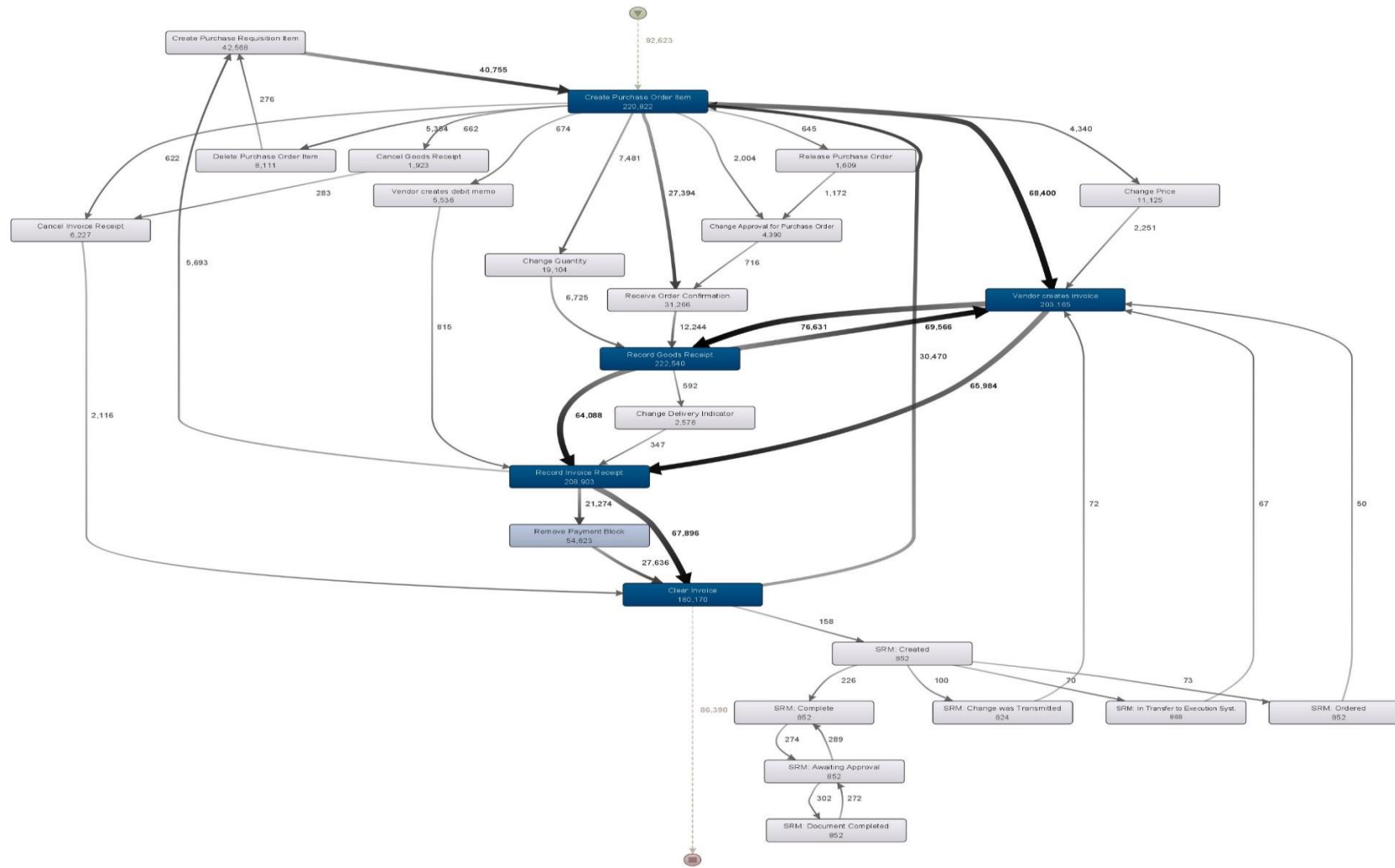
## Case item category – 3-way match invoice before

**(Values)** Purchase Items not requires goods receipt message/GR-based invoicing

**(Implication)** GR-based IV flag set to false and the Goods Receipt flags set to true.

Invoices entered before the goods are receipt(blocked until goods received

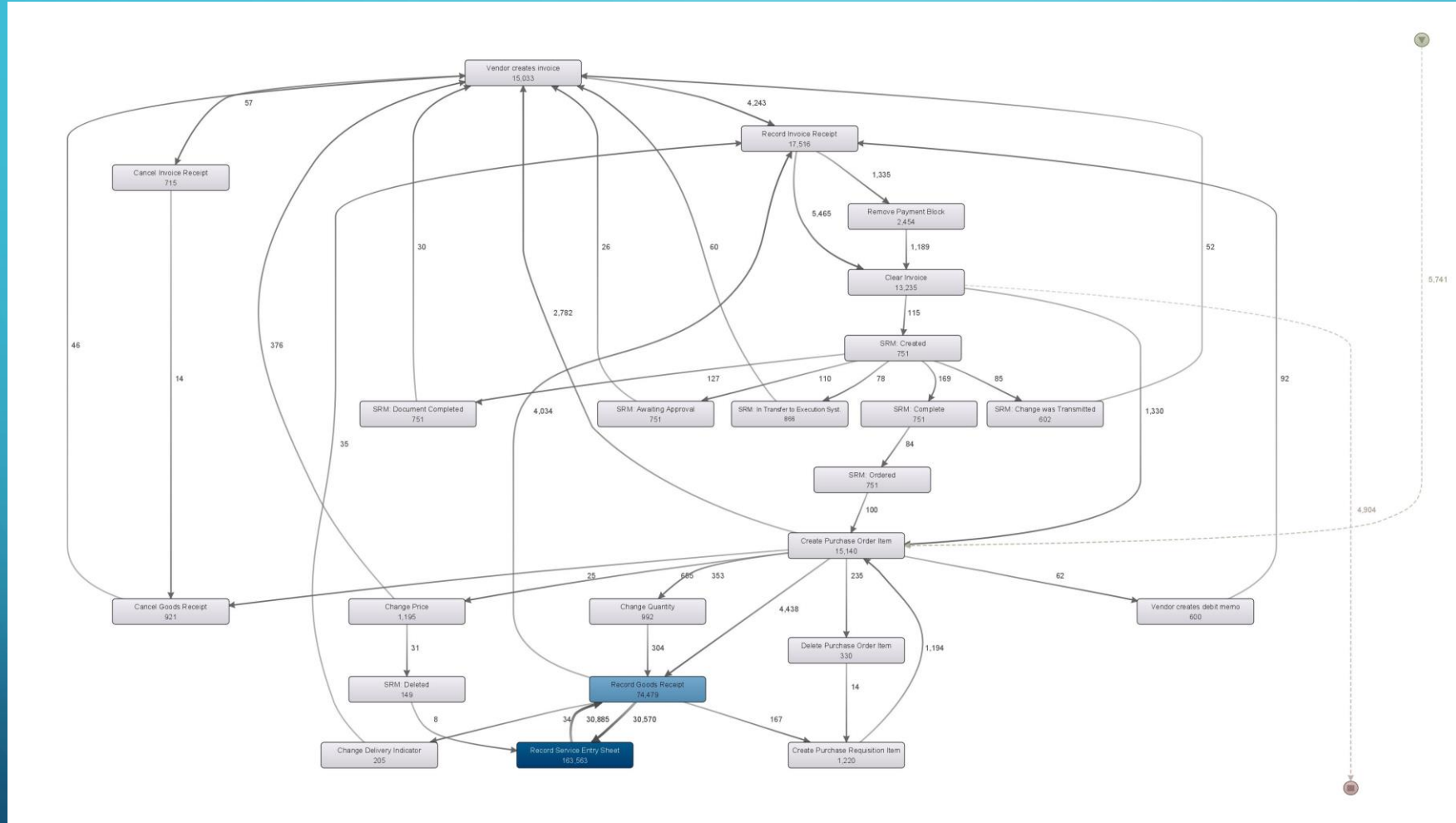
Unblocking is done by a user/a batch process at regular intervals



### Case item category – 3-way match, invoice after

**(Values)** An invoice receipt message ----matches ----- goods receipt message

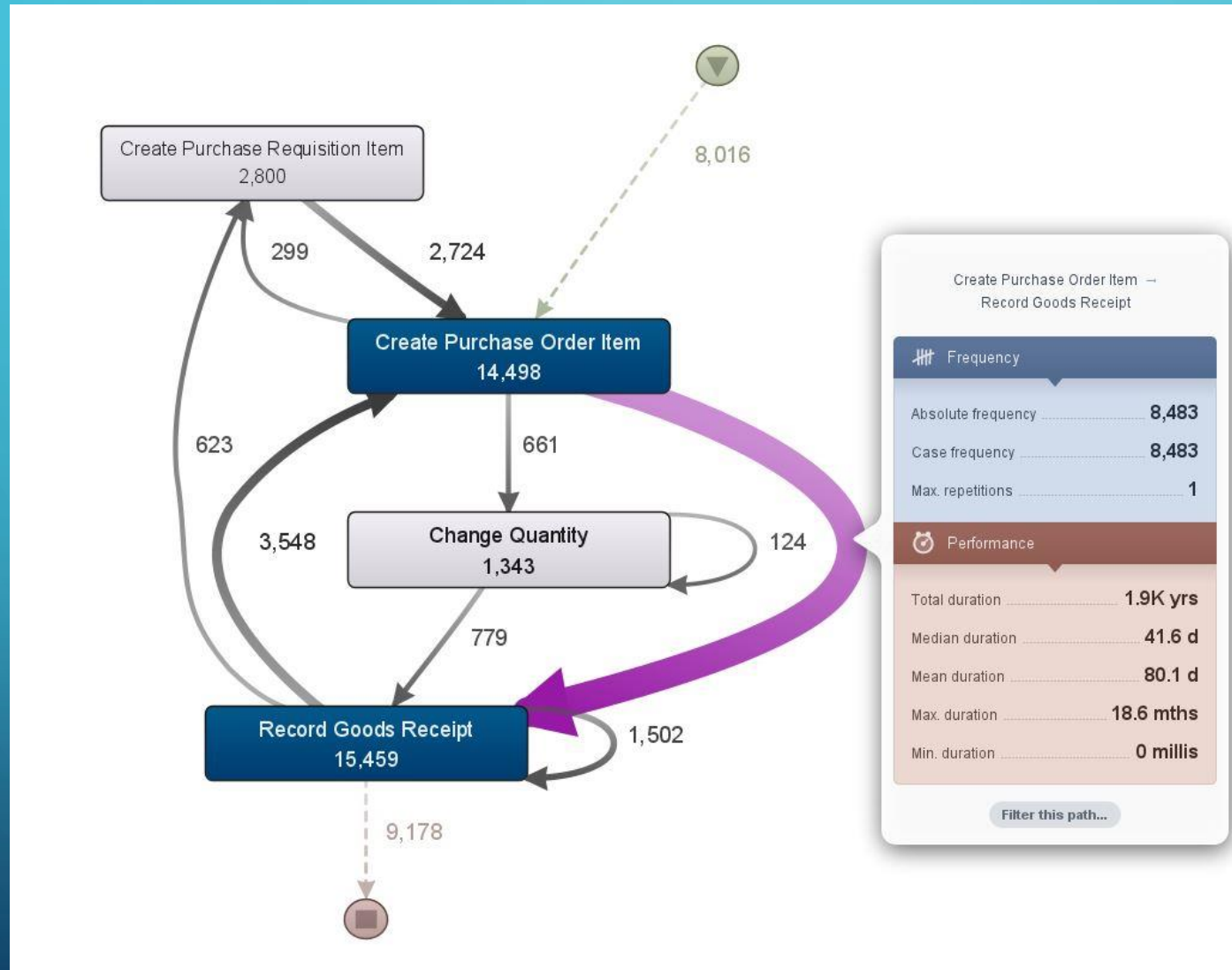
**(Implication)** GR-based flag and the Goods Receipt flags set to true



## Case Item Type - Consignment

(Values) classification text is PR

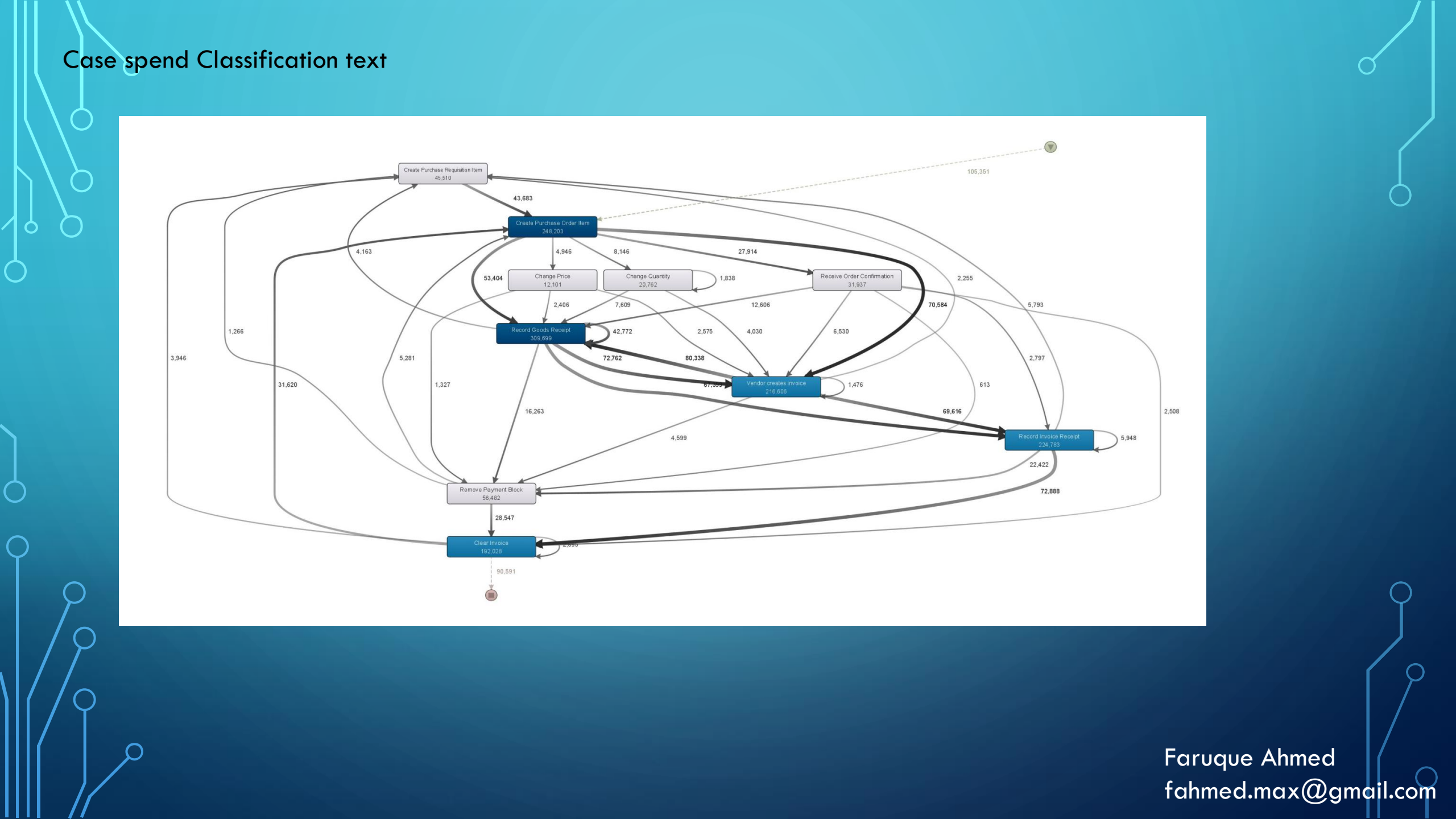
(Implication) GR indicator is set to False



## Case spend Classification text

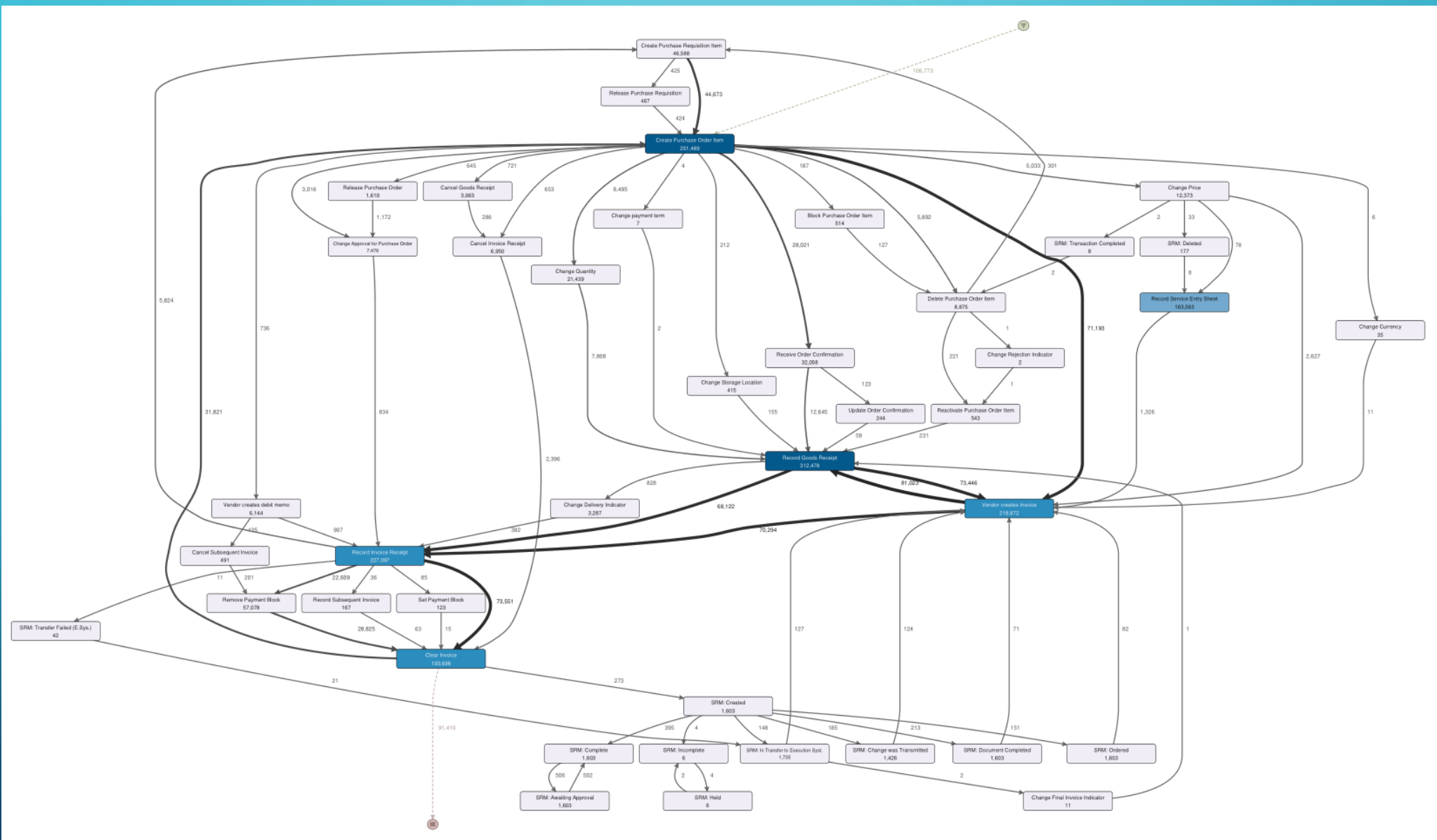
```
graph TD; Start(( )) --> C1[Create Purchase Requisition Item  
45,510]; C1 --> C2[Create Purchase Order Item  
248,203]; C2 --> C3[Change Price  
12,101]; C2 --> C4[Change Quantity  
20,762]; C3 --> C5[Record Goods Receipt  
309,699]; C4 --> C5; C5 --> C6[Vendor creates invoice  
216,696]; C5 --> C7[Record Invoice Receipt  
224,793]; C6 --> C7; C7 --> C8[Remove Payment Block  
56,482]; C8 --> C9[Clear Invoice  
192,028]; C9 --> End((End)); C1 --> C2; C2 --> C1; C2 --> C3; C3 --> C2; C2 --> C4; C4 --> C2; C3 --> C5; C5 --> C3; C4 --> C5; C5 --> C4; C5 --> C6; C6 --> C5; C5 --> C7; C7 --> C5; C6 --> C7; C7 --> C6; C7 --> C8; C8 --> C7; C8 --> C9; C9 --> C8; C9 --> End; C1 --> End; C2 --> End; C3 --> End; C4 --> End; C5 --> End; C6 --> End; C7 --> End; C8 --> End; C9 --> End;
```

Faruque Ahmed  
fahmed.max@gmail.com

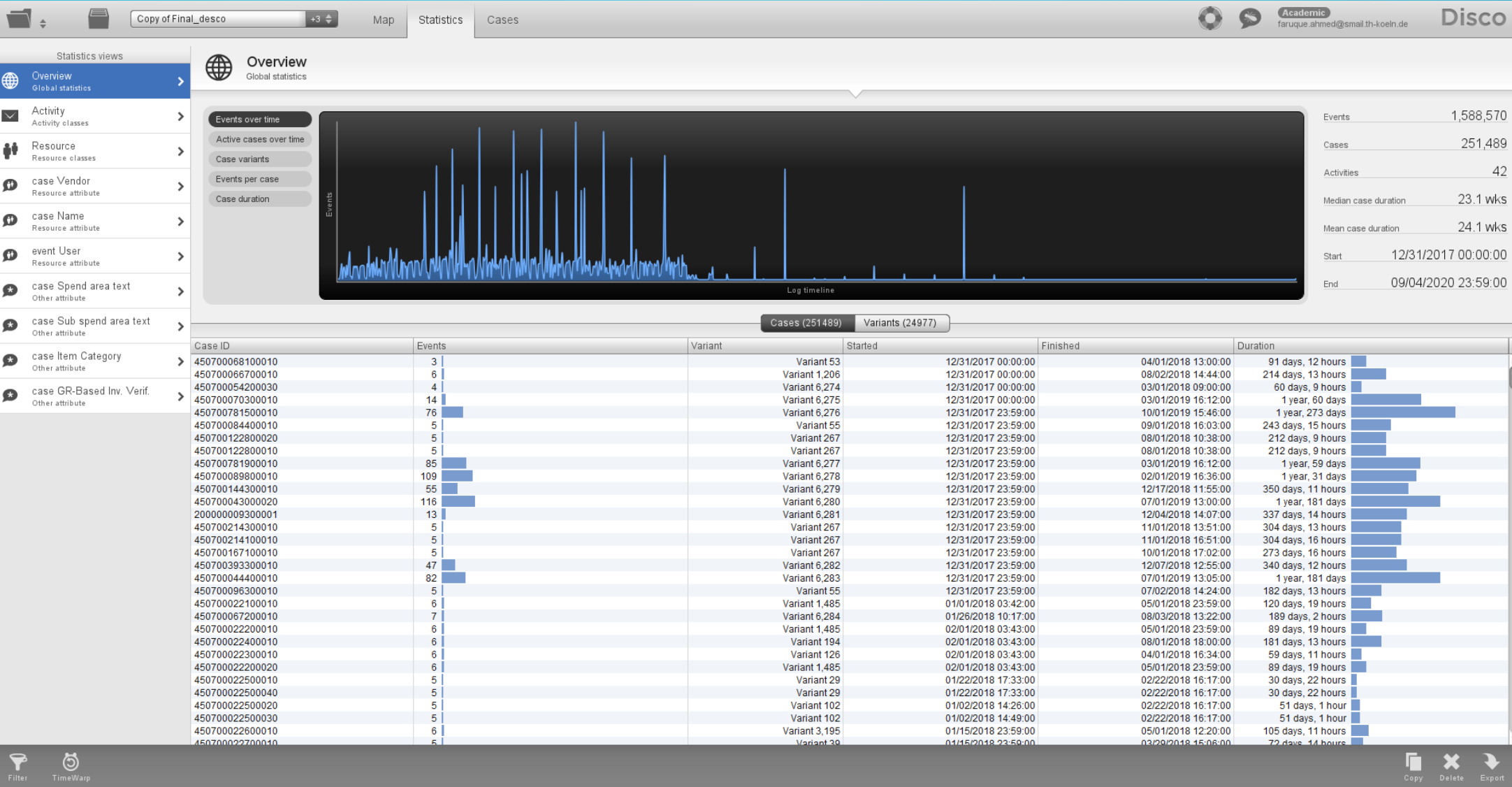




# Complete Process



Dashboard



[illegible]