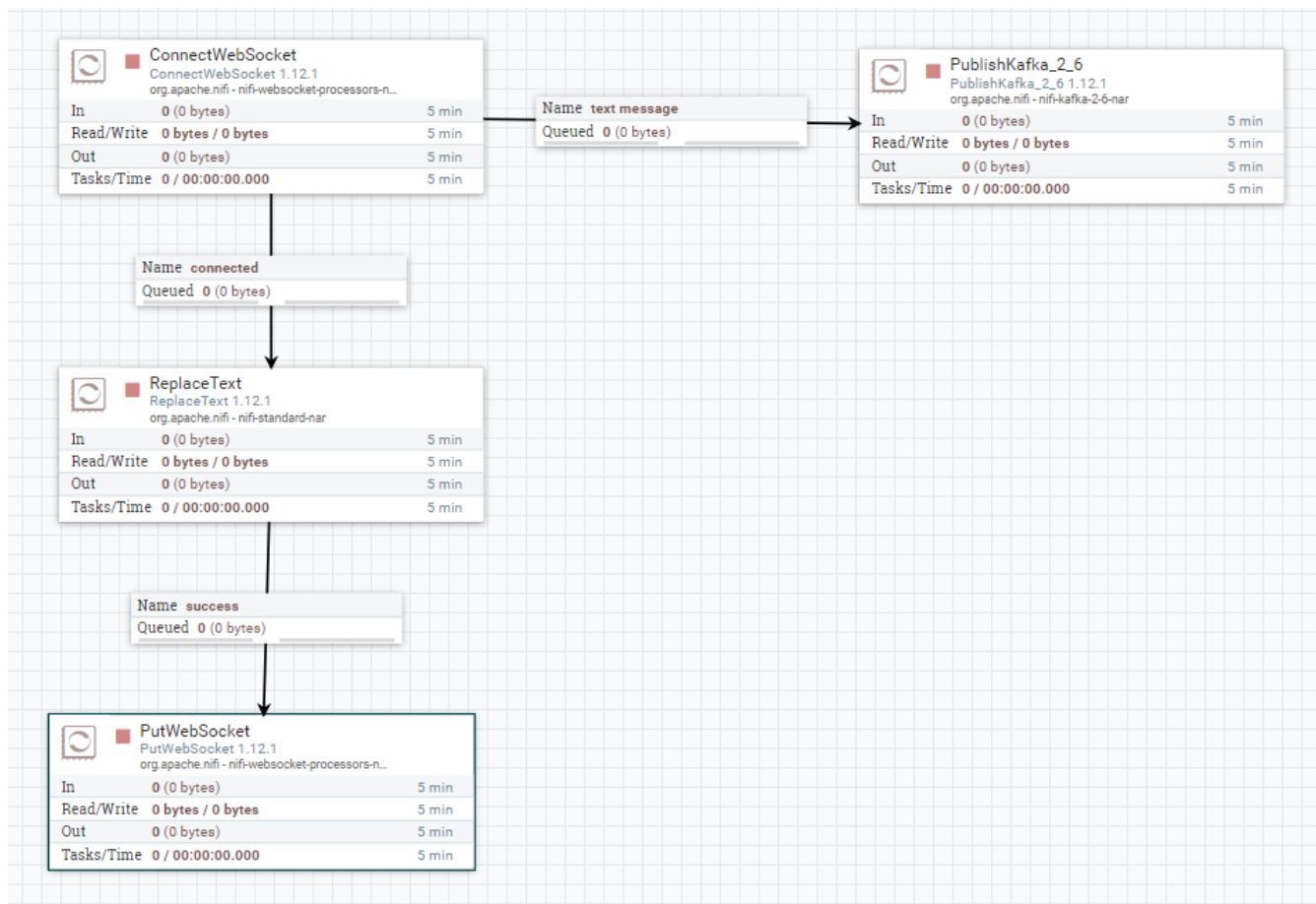# 1. NIFI flow

## 1.1. Import template
"sandbox_repo\training\GLBig_Data_ProCamp\home_work\HW2\nifi_kafka_hw2_template.xml". Create flow with imported template.



## 1.2. Create NIFI variable "CurrencyPair" with value "btcusd"



## 1.3. Configure and enable controller services: "JetyWebSocketClient", "StandardRegisteredSSLContextService"

## 2. Scripts preparation on GCP data proc node.

2.1. Copy shell and python scripts in home directory on GCP main node and grant permission:

Git Folder: "sandbox_repo\training\GLBig_Data_ProCamp\home_work\HW2"

Scripts:

- "btcusd_consumer_confl.py"
- "btcusd_consumer_confl.sh"
- "create_topic_btcusd.sh"

Execute command on main node (to grant needed permissions): "sudo chmod u+x *.sh"

# 3. Create Kafka topic.

Launch shell script "create_topic_btcusd.sh" to create topic "gcp.orders.fct.btcusd.0" and add current and "nifi" users to group "kafka":

```
drwxr-xr-x 3 mo_tarabanovskyi_gmail_com 2054902092  4096 Nov 28 10:16 snap
mo_tarabanovskyi_gmail_com@procamp-cluster-m:~$ ./create_topic_btcusd.sh
WARNING: Due to limitations in metric names, topics with a period ('.') or underscore ('_') could collide. To avoid
 issues it is best to use either, but not both.
Error while executing topic command : Topic 'gcp.orders.fct.btcusd.0' already exists.
[2020-12-13 01:20:35,302] ERROR org.apache.kafka.common.errors.TopicExistsException: Topic 'gcp.orders.fct.btcusd.0
' already exists.
 (kafka.admin.TopicCommand$)
check topic gcp.orders.fct.btcusd.0
gcp.orders.fct.btcusd.0
check nifi groups
nifi : nifi kafka
```

# 4. Install python modules "confluent-kafka", "pandas".

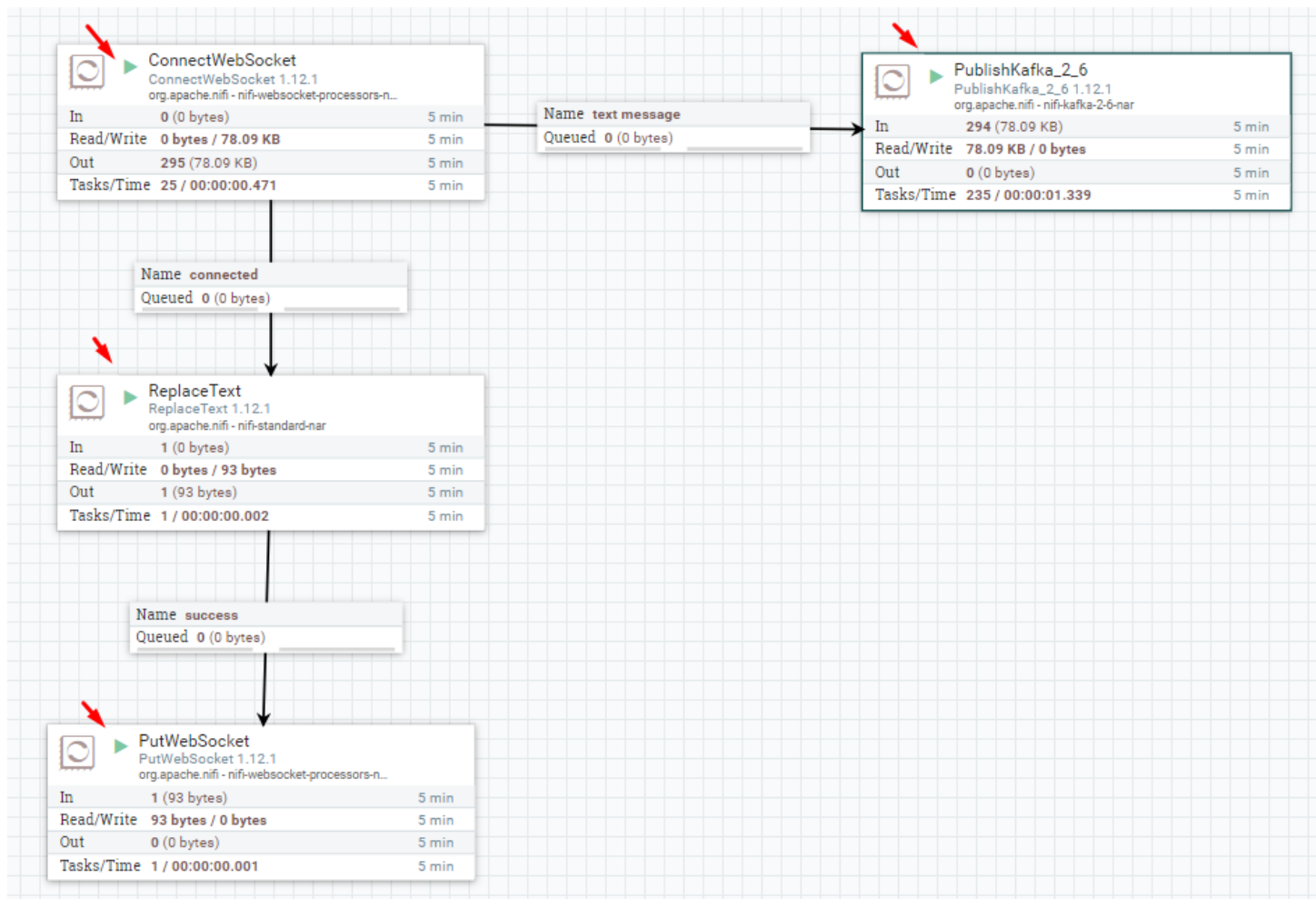4.1. Install python module "kafka-confluence":

```
sudo /opt/conda/default/bin/python -m pip install confluent-kafka
```

4.2. Reinstall python module pandas (to install actual version)

```
sudo /opt/conda/default/bin/python -m pip uninstall pandas
sudo /opt/conda/default/bin/python -m pip install pandas
```

# 5. Start kafka producer

Start all Apache NIFI flow processors:

## 6. Launch kafka consumer on GCP node.

Execute shell script "`btcusd_consumer_confl.sh`" on GCP master node.

```
mo_tarabanovskyi_gmail_com@procamp-cluster-m:~$ ./btcusd_consumer_confl.sh
2020-12-13 01:39:00,390 - INFO - ******* Set local kafka broker

****************** 1 messages are processed. Top 10 bitcoin transactions based on price field (descending):
                event data.id_str data.amount_str data.price_str
0  bts:subscription_succeeded        NaN             NaN            NaN


****************** 2 messages are processed. Top 10 bitcoin transactions based on price field (descending):
                event      data.id_str data.amount_str data.price_str
0         order_created  1306577185161216      0.15000000        18781.16
0  bts:subscription_succeeded          NaN             NaN            NaN


****************** 3 messages are processed. Top 10 bitcoin transactions based on price field (descending):
                event      data.id_str data.amount_str data.price_str
0         order_created  1306577185161216      0.15000000        18781.16
0         order_deleted  1306576997519360      0.14690970        18697.99
0  bts:subscription_succeeded          NaN             NaN            NaN


****************** 4 messages are processed. Top 10 bitcoin transactions based on price field (descending):
                event      data.id_str data.amount_str data.price_str
0         order_created  1306577186816000      0.53096254        18800.89
0         order_created  1306577185161216      0.15000000        18781.16
0         order_deleted  1306576997519360      0.14690970        18697.99
0  bts:subscription_succeeded          NaN             NaN            NaN


****************** 5 messages are processed. Top 10 bitcoin transactions based on price field (descending):
                event      data.id_str data.amount_str data.price_str
0         order_created  1306577186816000      0.53096254        18800.89
0         order_created  1306577185161216      0.15000000        18781.16
0         order_created  1306577186828288      0.47910926        18750.85
0         order_deleted  1306576997519360      0.14690970        18697.99
0  bts:subscription_succeeded          NaN             NaN            NaN
```

```
0  order_deleted  1306577430269955      8.97000000        18881.43


****************** 1044 messages are processed. Top 10 bitcoin transactions based on price field (descending):
        event      data.id_str data.amount_str data.price_str
0  order_created  1306577441411074      7.42000000        19148.64
0  order_deleted  1306574556114944      0.30000000        19010.68
0  order_created  1306577431527425      0.12700000        18980.00
0  order_deleted  1306577431527425      0.12700000        18980.00
0  order_created  1306577434177537      0.45000000        18889.50
0  order_deleted  1306577434177537      0.45000000        18889.50
0  order_created  1306577454764032      0.45000000        18889.50
0  order_created  1306577435574275      8.97000000        18887.54
0  order_created  1306577430269955      8.97000000        18881.43
0  order_deleted  1306577430269955      8.97000000        18881.43


****************** 1045 messages are processed. Top 10 bitcoin transactions based on price field (descending):
        event      data.id_str data.amount_str data.price_str
0  order_created  1306577441411074      7.42000000        19148.64
0  order_deleted  1306574556114944      0.30000000        19010.68
0  order_created  1306577431527425      0.12700000        18980.00
0  order_deleted  1306577431527425      0.12700000        18980.00
0  order_created  1306577434177537      0.45000000        18889.50
0  order_deleted  1306577434177537      0.45000000        18889.50
0  order_created  1306577454764032      0.45000000        18889.50
0  order_created  1306577435574275      8.97000000        18887.54
0  order_created  1306577430269955      8.97000000        18881.43
0  order_deleted  1306577430269955      8.97000000        18881.43
```

Note: if something went wrang, launch script with option "-v" for verbose output:

```
mo_tarabanovskyi_gmail_com@procamp-cluster-m:~$ ./btcusd_consumer_confl.sh -v
2020-12-13 01:43:39,338 - INFO - ******* Debug output is enabled
2020-12-13 01:43:39,338 - INFO - ******* Set local kafka broker
```

# 7. Launch kafka consumer on local PC. Windows OS example.

## (optional step, can be useful for debug).

7.1.

7.2. Install python modules on local-PC:

   "pandas", "kafka-consumer".

*"kafka-consumer" can be installed on python 3.7 from "whl" file*
*"sandbox_repo\training\GLBig_Data_ProCamp\infra\confluent\confluent_kafka-1.4.1-cp37-cp37m-*
*win_amd64.whl" (modules for version < 3.7 can be downloaded from confluent site).*

7.3. Up open VPN server on GCP master node

7.3.1.   Copy script on master node in home directory
      "sandbox_repo\training\GLBig_Data_ProCamp\infra\vpn\openvpn-configuration.sh"

7.3.2.   Launch script "sudo openvpn-configuration.sh" and fill parameters by default values (except: protocol - choose "TCP" and set client name)

```
Welcome to this OpenVPN road warrior installer!

Which IPv4 address should be used?
     1) 10.142.0.5
     2) 172.27.224.1
     3) 172.27.226.1
     4) 172.27.228.1
     5) 172.27.230.1
     6) 172.27.232.1
     7) 172.27.234.1
     8) 172.27.236.1
     9) 172.27.238.1
IPv4 address [1]:

This server is behind NAT. What is the public IPv4 address or hostname?
Public IPv4 address / hostname [34.73.161.236]:

Which protocol should OpenVPN use?
   1) UDP (recommended)
   2) TCP
Protocol [1]: 2

What port should OpenVPN listen to?
Port [1194]:

Select a DNS server for the clients:
   1) Current system resolvers
   2) Google
   3) 1.1.1.1
   4) OpenDNS
   5) Quad9
   6) AdGuard
DNS server [1]:

Enter a name for the first client:
Name [client]: alex

OpenVPN installation is ready to begin.
Press any key to continue...
Hit:1 http://us-east1.gce.archive.ubuntu.com/ubuntu bionic InRelease
Hit:2 http://us-east1.gce.archive.ubuntu.com/ubuntu bionic-updates InRelease
Hit:3 http://us-east1.gce.archive.ubuntu.com/ubuntu bionic-backports InRelease
Hit:4 https://storage.googleapis.com/goog-dataproc-bigton-repo-us-east1/1_5_deb10_20201018_013600-RC01_dataproc InR
```

7.3.3.   Copy client configuration on local PC

```
.............................................................................................................
.............................................................................................................
........................+++++
writing new private key to '/etc/openvpn/server/easy-rsa/pki/easy-rsa-32697.NCXVqt/tmp.5ydWqh'
-----
Using configuration from /etc/openvpn/server/easy-rsa/pki/easy-rsa-32697.NCXVqt/tmp.sErTm3
Check that the request matches the signature
Signature ok
The Subject's Distinguished Name is as follows
commonName            :ASN.1 12:'server'
Certificate is to be certified until Nov 26 21:20:52 2030 GMT (3650 days)

Write out database with 1 new entries
Data Base Updated

Using SSL: openssl OpenSSL 1.1.1  11 Sep 2018
Generating a RSA private key
......+++++
.................................................+++++
writing new private key to '/etc/openvpn/server/easy-rsa/pki/easy-rsa-304.SYCkbD/tmp.MaG6Fh'
-----
Using configuration from /etc/openvpn/server/easy-rsa/pki/easy-rsa-304.SYCkbD/tmp.cwOk3a
Check that the request matches the signature
Signature ok
The Subject's Distinguished Name is as follows
commonName            :ASN.1 12:'alex'
Certificate is to be certified until Nov 26 21:20:52 2030 GMT (3650 days)

Write out database with 1 new entries
Data Base Updated

Using SSL: openssl OpenSSL 1.1.1  11 Sep 2018
Using configuration from /etc/openvpn/server/easy-rsa/pki/easy-rsa-361.DVw8Qc/tmp.5wxhBe

An updated CRL has been created.
CRL file: /etc/openvpn/server/easy-rsa/pki/crl.pem


Created symlink /etc/systemd/system/multi-user.target.wants/openvpn-iptables.service → /etc/systemd/system/openvpn-
iptables.service.
Created symlink /etc/systemd/system/multi-user.target.wants/openvpn-server@server.service → /lib/systemd/system/ope
nvpn-server@.service.

Finished!

The client configuration is available in: /home/mo_tarabanovskyi_gmail_com/alex.ovpn
New clients can be added by running this script again.
```

7.3.4.    Change vpn server configuration and reboot master node

Server configuration example: "sandbox_repo\training\GLBig_Data_ProCamp\infra\vpn\server.conf"

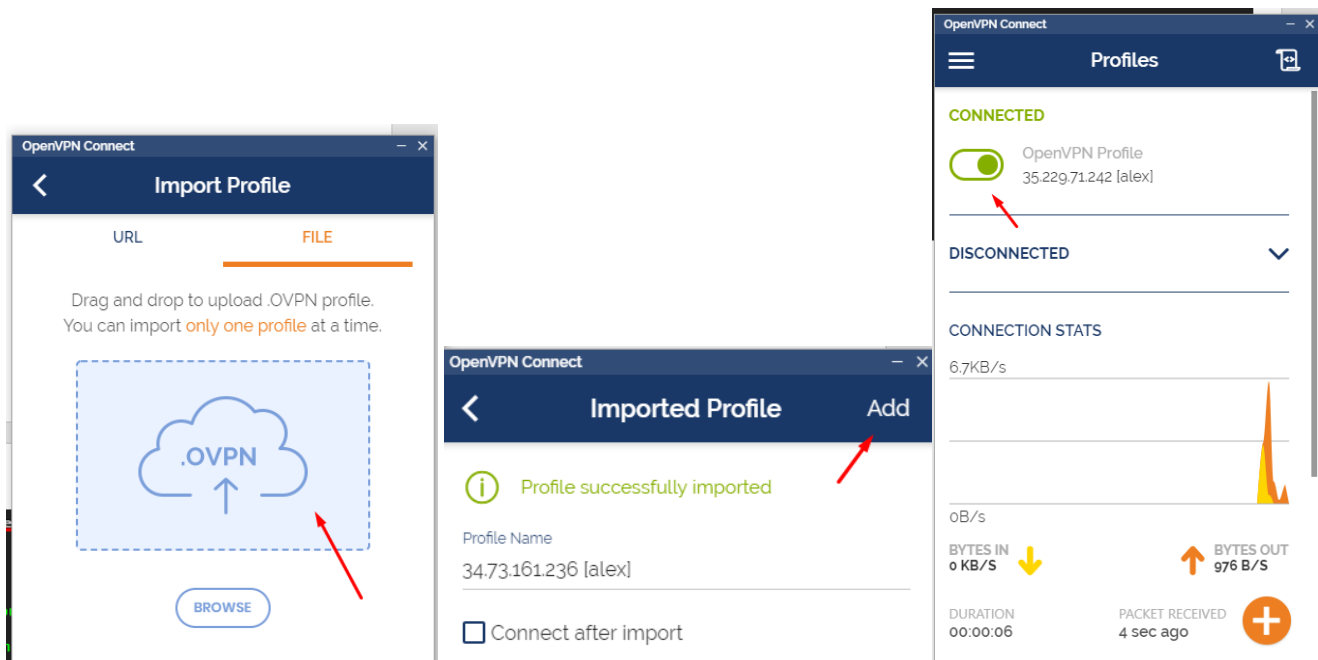The configuration location is "`/etc/openvpn/server/server.conf`" on GCP master node.

```
mo_tarabanovskyi_gmail_com@procamp-cluster-m:~$ cat /etc/openvpn/server/server.conf
local 10.142.0.5
port 1194
proto tcp
dev tun
ca ca.crt
cert server.crt
key server.key
dh dh.pem
auth SHA512
tls-crypt tc.key
topology subnet
server 10.8.0.0 255.255.255.0
push "route 10.142.0.0 255.255.255.0"
#push "redirect-gateway def1 bypass-dhcp"
ifconfig-pool-persist ipp.txt
push "dhcp-option DNS 169.254.169.254"
push "route 169.254.169.254"
keepalive 10 120
cipher AES-256-CBC
user nobody
group nogroup
persist-key
persist-tun
status openvpn-status.log
verb 3
crl-verify crl.pem
```

1 .Should be added for GCP nodes network routing. Check subnet addresses

2. Should be commented to avoid Internet traffic routing via VPN

3. Should be added for GCP nodes names routing ("*.internal" host names)

7.4. Install Open VPN client and import client configuration file (described in 7.2.3) and enable connection



7.5. Change GCP master node name in python script
"sandbox_repo\training\GLBig_Data_ProCamp\home_work\HW2\btcusd_consumer_confl.py":

```
36      if args.remote_bootstrap:
37          logger.info('****** Set remote kafka broker')
38          # Should be set master node internal hostname in format <host name>.<region>.<project id>.internal
39          # open vpn server should use configuration:
40          #    sandbox_repo\training\GLBig_Data_ProCamp\infra\vpn\openvpn-configuration.sh
41          lv_bootstrap_servers = 'procamp-cluster-m.us-east1-b.c.bigdata-procamp-1add8fad.internal'
42      else:
43          logger.info('****** Set local kafka broker')
44          lv_bootstrap_servers = 'localhost:9092'
45
```

7.6. Launch consumer with option "-r" (remote broker) –
"sandbox_repo\training\GLBig_Data_ProCamp\home_work\HW2\btcusd_consumer_confl.py"