

We use Equations 3.1-3.4 to moderate the network congestion in the following scenarios:

- A. Normal case (*many edge devices are running with fluctuating degrees of load*): in this scenario, from the equations cited above, some of the edge devices will favour local processing more while others will favour cloud processing as shown in Fig. 3. Therefore, no congestion is anticipated, as these events are random.
- B. Light Load case (*many edge devices are idle and not streaming most of the time as is the case with event-driven embedded devices*): in this situation, the system works well without network congestion as cloud resources will be required only by busy edge devices as we anticipate that the occurrence of these events is random.
- C. Worst case (*many edge devices are working and need access to the cloud resources because they have reached the limit of the resources dedicated for anomaly detection in the edge device*): this represents a burst scenario and we anticipated that while this is rare, it is also possible. Therefore, the architecture uses the  $TDD_1$  of Equation 3.4 to handle such situation. If  $TDD_1$  becomes bigger than the permissible latency, that particular detection operation is dropped silently while the model uses the Equations 3.3 and 3.4 to balance the network.