

Distributed Order Dependency

Sebastian Schmidl, Juliane Waack¹

Abstract: Abstract goes here.

Keywords: Actor Model; Akka; Distributed Computing; Parallelization

1 Introduction

1.1 Order Dependencies

- What are order dependencies (ODs)?
- Why are they useful?
 - Query optimization
 - Data quality (find constraints, help understand semantics)
 - support index selection

1.2 Order Dependency Discovery

What are the challenges?

- Defining Minimality (Unhelpful definition of minimality in Naumann paper, incorrect one in Consonni paper)
- Large searchspace (Long runtime and large memory-needs)
- current algorithms slow on large dataset
- Therefore need scalability -> distribution

1.3 Distribution

- What are distributed Systems? Why are they useful?

¹ Hasso-Plattner-Institut, University of Potsdam, Prof.-Dr.-Helmert-Str. 2-3, 14482 Potsdam, {sebastian.schmidl, juliane.waack}@student.hpi.de

- What is the actor model? Why use it here?

2 Related Work

In comparison to the related field of functional dependencies, literature on ODs is comparatively rare. The idea of considering the ordering of attributes in a relation as a kind of dependency was first introduced in 1982 by Ginsburg; Hull [GH83]. Their work formalizes the so-called *point-wise ordering* with a complete set of inference rules and shows that the problem of inference in their formalism is co-NP-complete [GH83]. Ginsburg; Hull's definition of *point-wise ordering* specifies that a set of attributes in a relation orders another set of attributes.

In 2012, Szlichta et al. [SGG12] introduced a definition for ODs, which considers lists of attributes instead of sets of attributes. This leads to a lexicographical ordering of tuples as generated by the `ORDER BY` operator in SQL. The definition of ODs by Szlichta et al. was used throughout the following research in this area [Co19; LN16; Sz17] and is also the basis of this work.

The problem of efficient discovery of ODs in existing datasets using Szlichta et al.'s definition was first approached by Langer; Naumann [LN16]. Their algorithm is called `ORDER` and first computes all potential OD candidates and then traverses the candidate lattice in a bottom-up manner employing pruning rules. `ORDER` has a factorial worst-case complexity over the number of attributes in the relation. Unfortunately, the aggressive pruning rules of `ORDER` affect the completeness of their results [Sz17].

Szlichta et al. present another approach to OD discovery called `FASTOD`. It is complete and faster than `ORDER`. `FASTOD` uses a polynomial mapping of list-based ODs to a set-based canonical form, which reduces the algorithm's worst-case complexity to $O(2^{|n|})$, in the number of attributes n .

Most recently, Consonni et al. [Co19] presented a third approach on OD discovery using the concept of order compatibility dependencies (OCDs), called `OCDDISCOVER`. It has a factorial worst-case runtime, but Consonni et al. showed that it can outperform `FASTOD` on some datasets. `OCDDISCOVER` is capable of running multi-threaded, which makes it easier to adapt the algorithm to a distributed system. This was the state of published research when we started on our project, which is why we based our own algorithm on `OCDDISCOVER`. Since then, Godfrey et al. [Go19] have published an Errata Note demonstrating an error in Consonni et al.'s proof of minimality, which renders the results of `OCDDISCOVER` incomplete. We still base our approach on `OCDDISCOVER`, since the goal of our work is the distribution of an already existing algorithm with a focus on reliability and scalability.

3 Approach

3.1 OCDDISCOVER

The OCDDISCOVER algorithm creates the search space for finding ODs over OCDs. This can be done, because they prove that whenever there is an OD, there is also an OCD and a functional dependency. Using OCDs reduces the search space, because they are symmetrical.

Consonni et al. build the search tree, by testing each candidate for being order compatible. If it is not, no children are generated for this candidate. If it is order compatible the candidate is checked for being order dependent. If it is we have found a minimal OD and no children need to be checked, as they would not be minimal. If it is not, children to be checked later are generated. This happens by adding a column that does not yet appear on either side of the candidate once to its left and once to its right side [Co19].

Consonni et al. proved in their paper, that they would be able to find all minimal ODs using this method of generating candidates. However, as Godfrey et al. showed in their Errata Note, Consonni et al. made a mistake in one of their proofs and this method of building the search space does not create minimal ODs with repeated prefixes [Go19].

3.2 Our architecture

3.2.1 System architecture

The cluster we use has a *Peer-to-peer* architecture. Every node has the same structure and works on the same kinds of tasks. Every node also hold the complete dataset, because the checking of OD candidates requires access to the contents of all columns. The *Peer-to-peer* architecture allows every single node to fail without work getting lost or the algorithm needing to be restarted.

The only circumstance during which this is not the case is when the systems is newly started. At the beginning a single node, which we call *Seed node*, is responsible for loading the data, the initial pruning step and the generation of the first set of candidates. If the *Seed node* were to fail during these first steps, the cluster would need to be restarted. As they usually only take a few seconds, though, this seemed to be a reasonable constraint.

To begin working, every node needs the dataset, the information which columns can be pruned and a first set of candidates to check. All nodes that are not the *Seed node* can get the dataset from another node via streaming. Once it is done with pruning, the *Seed node* sends the pruning results to all other nodes in a broadcast. A node that joins later can ask any other for this information. The distribution of candidates from the *Seed node* to all other nodes in the cluster is done using *Work Stealing* (Sect. 3.3.1), which is described in more detail later.

3.2.2 Node Architecture

Every node in the cluster shares the same actors. These are both helpful for parallelizing the tasks each node work on as well as structuring the communication in between nodes.

Master The *Master actor* is responsible for holding the state of a node. It contains the *Candidate queue* and *Pending queue* and distributes the candidates to the *Worker actors* to be worked on in parallel. We use a pull-based approach for this distribution of tasks. This means that the *Master* only ever responds to *Workers* asking for candidates and therefore does not have to monitor the state of every single *Worker* of the node. The *Master* also deals with getting new candidates using *Work Stealing* (Sect. 3.3.1), when the *Candidate Queue* is empty and initiates the node shutdown when the *Downing Protocol* (Sect. 3.3.2) reveals that all ODs have been found.

Worker A node can have an arbitrary number of *Worker actors*, though the number of available threads - 1 is usually a good estimation for effective resources usage. The *Workers* check candidates, which they get from the *Master* and generate new ones, which they send back for the *Master* to append to the *Candidate Queue*. This checking and generating of candidates is the main task of each node and therefore requires the most resources. All ODs and OCDs the *Worker* finds are sent to the *Result Collector* to be recorded into a file.

Data Holder The *Data Holder actor* is responsible for loading and holding the dataset, for which the ODs are to be found. A reference to this dataset is shared with the *Workers* who need it to check candidates. If a path to the dataset is provided to the node on start up, it loads it from the file. If not, the *Data Holder* asks one of their neighbors in the cluster to provide them the dataset via streaming. Once the dataset is loaded, the *Data Holder* notifies the *Master* that the node is ready to start working on checking candidates.

Result Collector The *Result Collector actor* gets sent all found ODs and OCDs by the *Workers*. It collects and counts these results and writes them to a file so they can be recovered even if the node fails unexpectedly.

Cluster Listener The *Cluster Listener actor's* task is to determine and share who the neighbors of its node are. Whenever a new node joins the cluster all previous send it references to their *Master* and *State Replicator actors*. The new *Cluster Listener* sends back the corresponding information for its node.

When the ordering of nodes changes because a node joined or left the cluster, the updated neighbors are sent to the *State Replicator*. In contrast, the *Data Holder* only receives information about the neighboring nodes, when it specifically asks for them, because it only needs the information in the beginning to know whom to ask for the dataset.

State Replicator The *State Replicator* holds the state of its neighboring nodes and shares the *Master's* state with them. This enables us to recover work from unexpectedly failed nodes and is explained in more detail in the *State Replicaton Protocol* (see Sect. 3.3.3).

Reaper The *Reaper Actor* watches all other actors and cleanly shuts down the system once all of them are terminated.

3.3 Communication Protocols

3.3.1 Work Stealing Protocol

When a node is out of work, either because it just joined the cluster or because all the candidates it checked got pruned and did not generate any new candidates to check, it tries to take over some of the work of the other nodes. This process is called *Work Stealing*. We use it to ensure that no node is idle while the others are still checking candidates and to balance the workload over the cluster. The workload of a node is defined by the number of candidates waiting to be processed in its queue. In the desired outcome to add to its own queue of the *Work Stealing Protocol*, each node ends up with a similar workload. To achieve this goal the *Master Actors* of the different nodes communicate their workloads with each other.

The *Work Stealing Protocol* is initiated when a node's *Candidate Queue* is empty, even if some of its workers are still processing candidates and might return new ones soon. This node (the *Thief node*) sends a message to all other nodes, asking for the current size of their *Candidate Queues*. It also sets a timeout after which it processes the answers of the other nodes. During this processing step it calculates the average *Candidate Queue's* length of all the nodes that answered and itself. To improve the balancing of workloads over the cluster, the *Thief node* only takes candidates from nodes that have a *Candidate Queue* of above average length and only takes so many candidates that its own *Candidate Queue's* length does not exceed the average. To get the candidates from another node's *Candidate Queue*, the *Thief node* sends a message with the amount of candidates it wants to take to the node it wants to take them from. This node, the *Victim node*, then moves the asked for candidates from its *Candidate Queue* into its *Pending Queue* and sends them to the *Thief node*. However, since the *Victim node's* queue might have changed since it first sent its length to the *Thief node*, the *Victim node* sends back at most half of the candidates in its queue. Once the *Thief node* receives the candidates, it adds them to its own queue and sends a message acknowledging the receipt to the *Victim node*, which then deletes the candidates from its *Pending queue*. To ensure the send candidates arrive and will get processed, the *Victim node* watches the *Thief node*. In the case that the *Thief node* fails before it could acknowledge the receipt of the candidates, the *Victim node* moves them back to its own *Candidate queue* from its *Pending queue*.

3.3.2 Downing Protocol

When all ODs have been found and there are no more candidates to check, the system should shut itself down. A node starts the *Downing Protocol*, when it has no more candidates in its *Candidate Queue* and *Pending Queue* and an attempt to get more work using *Work Stealing* was unsuccessful. Checking the *Pending Queue* is important because candidates currently being processed could generate new candidates and therefore new work.

To ensure that there are no more candidates to be processed and no more candidates being generated in the whole cluster, the node has to ask every other node, whether they are also out of candidates. If another node still holds unprocessed candidates or is still in the process of creating new candidates, these could be redistributed over the cluster using *Work Stealing* so the node can continue to work instead of shutting itself down. The node waits for a response from every other node in the cluster, also being aware of nodes leaving the cluster and subsequently not awaiting their response anymore.

If any of the other nodes respond with a message saying that they are still holding or processing candidates, the node switches from the *Downing* to the *Work Stealing Protocol*. If all nodes have either left the cluster or responded that they have no more work, the node shuts itself down.

3.3.3 State Replication Protocol

To ensure any of the nodes of the cluster can fail without any work getting lost, every node regularly replicates its state to two other nodes. Because of this, three specific nodes would have to fail simultaneously to permanently lose any unprocessed candidates. In this case, everything would need to be restarted to ensure complete results.

Which nodes the state is sent to is decided by their addresses. All nodes in a cluster are ordered by their addresses and every node replicates its state to both its neighbors in this ordering. The first and last node in this ordering are also each others neighbors.

Every node has a *State Replication Actor* responsible for regularly sending the node's state to its neighbors and processing these neighbor's states when they are sent. It gets the information of who its neighbors are from the *Cluster Listener Actor*.

In regular intervals, the *State Replicator* asks the *Master Actor* for its current state. The state is a combination of its *Candidate Queue* and *Pending Queue*. On receiving the *Master's* current state, the *State Replicator* sends to both its neighbors with a version number that is continually increased.

When a *State Replicator* receives the state of one of its neighbors, it checks whether the version number in the new message is higher than the version number of the last state it received from that actor. If this is the case, it saves the new state and version number. If not, it discards the new message, because it already has a more up to date state for that actor.

A node joins the cluster When a new node joins the cluster, the *Cluster Listener* updates its ordering of all the nodes in the cluster. It then sends the *State Replicator* the updated neighbor information. If the nodes neighbor has changed, the *State Replicator* removes the state of the node that is no longer its neighbor and sends its own state to its new neighbor.

A node leaves the cluster When a node leaves the cluster, its two neighbors have to decide which of them will take on the work the leaving node did not finish. They are alerted to their neighbor's leaving by the *Cluster Listener* who also immediately tells them, who their new neighbor is. This new neighbor was the leaving nodes other neighbor. To determine which of the two has the most up to date state of the leaving node, they share the version number of the leaving node's state. The node with the higher version number adds the leaving nodes state to its own *Candidate Queue* to ensure all candidates get processed. The node with the lower version number deletes the state of the leaving node. The two nodes then share their state with each other.

4 Evaluation

In this section, we present an initial evaluation of the performance and fault-tolerance of our algorithm. In Sect. 4.1, we evaluate the correctness of the results produced by our algorithm when run in different settings and under faults, Sect. 4.2 then compares the performance with our baseline OCDDISCOVER by Consonni et al. Sect. 4.3 explores the scalability of DODO across multiple nodes. The code for OCDDISCOVER was provided to us by the authors.

To evaluate the performance of the different OD discovery approaches, we measure the overall runtime the algorithm needed to compute all ODs and OCDs. The runtime of the algorithm includes startup and shutdown times and the algorithm must terminate to produce a valid result. As a consequence, when running DODO in a distributed setting, this measurement will include node discovery and shutdown synchronization times.

All experiments were performed on homogeneous nodes with 20 cores and 64 GB of main memory. DODO is implemented in Scala and the nodes run it with Java 1.8 with the Java JVM heap space limited to 40 GB. We adopted the *flight_1k* dataset available from the Information Systems Group at Hasso Plattner Institute³ to reduce the computation time to an order of minutes. This was achieved by removing sparse columns⁴. We ended up with a dataset consisting of 1000 rows and 36 columns: *flight_r1k_c36*.

³ <https://hpi.de/naumann/projects/repeatability/data-profiling/fds.html>

⁴ Removed columns: ArrivalDelayGroups, CancellationCode, Cancelled, CarrierDelay, CRSArrTime, DepartureDelayGroups, DepTimeBlk, Dest, DestAirportID, DestAirportSeqID, DestCityMarketID, DestCityName, DestState, DestStateFips, DestStateName, DestWac, DistanceGroup, Div1Airport, Div1AirportID, Div1AirportSeqID, Div1LongestGTime, Div1TailNum, Div1TotalGTime, Div1WheelsOff, Div1WheelsOn, Div2Airport, Div2AirportID, Div2AirportSeqID, Div2LongestGTime, Div2TailNum, Div2TotalGTime, Div2WheelsOff, Div2WheelsOn, Div3Airport, Div3AirportID, Div3AirportSeqID, Div3LongestGTime, Div3TailNum,

4.1 Correctness and fault-tolerance

Because we based our approach on OCDDISCOVER by Consonni et al., we don't expect our algorithm to output a complete minimal set of ODs. Our goal is to output the same OCDs and ODs as OCDDISCOVER.

We find that DODO computes the same OCDs and ODs as OCDDISCOVER, independent of whether DODO is run on a single node or distributed across multiple nodes. We did not compare the actual results for our medium-size test dataset *flight_r1k_c36*, because it contains too many ODs. But Tab. 1 shows, that DODO finds the same number of ODs in the dataset as OCDDISCOVER. Column four in Tab. 1 indicates the sum of found OCDs and ODs. *iw, jn, kfailure(s)* indicates an experiment with DODO using *i* workers per node, the system consisting of *j* nodes overall and *k* nodes of them were killed during the experiment.

Experiment	CCs	OECs	OCDs + ODs	runtime
OCDDISCOVER, 8w	7	5	49 514	7 m 09 s
8w, 1n, 0failure	7	5	49 514	3 m 42 s
8w, 8n, 0failure	7	5	49 514	49 s
8w, 8n, 1failure	7	5	50 055	54 s
8w, 8n, 2failures	7	5	48 374	55 s
8w, 8n, 3failures	7	5	48 352	52 s

Tab. 1: Result comparison between OCDDISCOVER, DODO, and DODO in the presence of node failures.

If DODO experiences node failures, it redistributes the OCD candidates to the remaining nodes using the state replication protocol (see Sect. 3.3.3). This leads to duplicate checks and along with that to the same OD being found multiple times, as is the case in the experiment 8w, 8n, 1failure. DODO finds more ODs in this experiment as there are in the dataset. This is due to the fact that we do not yet perform deduplication of the found ODs.

The experiments with two and more failing nodes do not find all ODs in the dataset. Unfortunately, we performed the experiments manually and the first node failure was introduced before the node was able to replicate its state to the other nodes in the cluster. This caused the system to loose all OCD candidates of this node. But DODO was able to recover from the following node failures, because those nodes were already able to replicate their state.

Div3TotalGTime, Div3WheelsOff, Div3WheelsOn, Div4Airport, Div4AirportID, Div4AirportSeqID, Div4LongestGTime, Div4TailNum, Div4TotalGTime, Div4WheelsOff, Div4WheelsOn, Div5Airport, Div5AirportID, Div5AirportSeqID, Div5LongestGTime, Div5TailNum, Div5TotalGTime, Div5WheelsOff, Div5WheelsOn, DivActualElapsedTime, DivAirportLandings, DivArrDelay, DivDistance, DivReachedDest, FirstDepTime, LateAircraftDelay, LongestAddGTime, NASDelay, SecurityDelay, TaxiIn, TaxiOut, TotalAddGTime, WeatherDelay, WheelsOff, WheelsOn

In Tab. 1, we already notice that DODO in a single nodes setup with eight workers is nearly 2x as fast as OCDDISCOVER with eight threads. Running DODO on multiple nodes decreases the time even more as one would expect. If there are no failures, DODO outputs the correct results independent of the number of nodes it is deployed on.

4.2 Comparing with OCDDISCOVER

In Fig. 1, we compare the performance of DODO with OCDDISCOVER on a single node using our test dataset. OCDDISCOVER is only capable of running on a single-node setup, but can scale by increasing the number of used threads. We measure the runtime of the algorithms to find all ODs with a different degree of parallelism. OCDDISCOVER scales by increasing the number of threads and DODO scales with the number of workers in a single node setup.

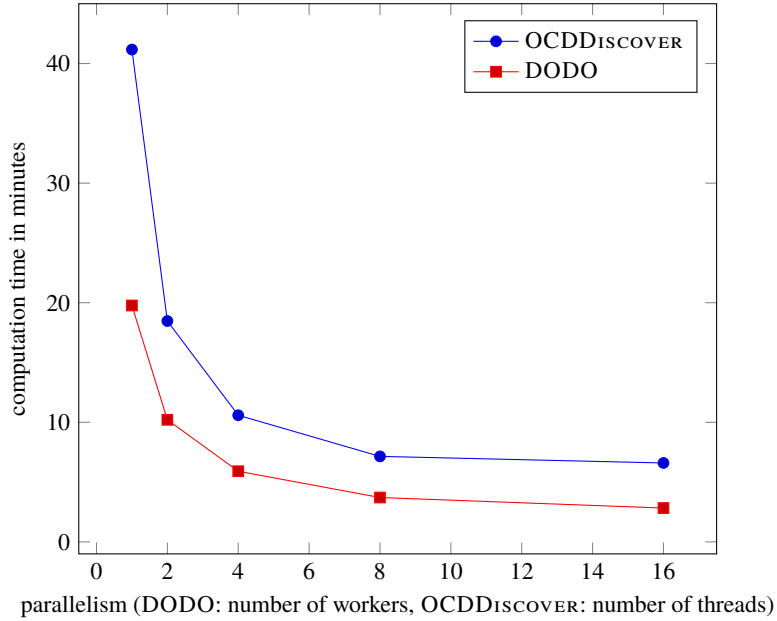


Fig. 1: Single node runtime comparison of our approach and the baseline algorithm OCDDISCOVER with the test dataset *flight_r1k_c36*.

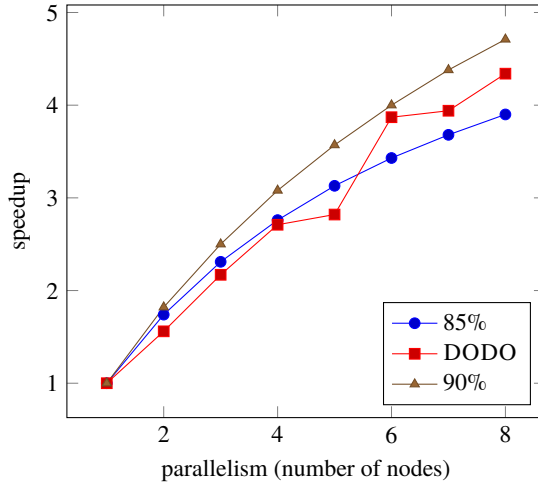
We find that DODO with a single worker already outperforms the single-threaded version of OCDDISCOVER by over a factor of two. OCDDISCOVER switches to a completely different implementation when it is run multi-threaded. This version performs slightly better compared to DODO. Here, DODO can only achieve a performance increase by a factor of around 1.8 compared to OCDDISCOVER. But besides that, DODO scales better than OCDDISCOVER. DODO achieves a performance increase of a factor of 2.8 when

run with 18 workers compared to OCDDISCOVER with 18 threads. This shows that our implementation using the actor model with no explicit synchronization barriers is superior to the explicit multi-threading implementation by Consonni et al.

4.3 Scalability

In Sect. 4.2, we already showed that DODO outperforms our baseline OCDDISCOVER in a single node setup. In this section, we analyze the scalability of DODO when run on multiple nodes in a cluster setup.

Fig. 2 shows the results of our cluster scalability experiments using our test dataset *flight_r1k_c36*, where we measured the runtime of DODO with different cluster sizes. Fig. 2a plots the speed of running DODO on multiple nodes and Fig. 2b shows the actual runtime of DODO.



(a) Computation time speedup of DODO compared to theoretical 85% and 90% speedup.

# nodes	computation time
1	3 m 37 s
2	2 m 19 s
3	1 m 40 s
4	1 m 20 s
5	1 m 17 s
6	56 s
7	55 s
8	50 s

(b) Computation time for finding all ODs in the dataset of different cluster sizes.

Fig. 2: Scaling our algorithm over the number of nodes and keeping the number of workers fixed at eight workers per node using the test dataset *flight_r1k_c36*.

We find that using DODO in a cluster setup greatly reduces overall computation time despite the additional effort spend on node discovery, state replication and inter-node synchronization. Finding all ODs in our test dataset takes about three and a half minutes on a single node. This time is already reduced to under a minute when using six or more nodes in the DODO cluster. As we can see in Fig. 2a, DODO can achieve a speedup between 85% and 90%. The outliers in the diagram can be explained by our downing protocol implementation, which waits for all nodes to come to a mutual decision that all candidates were processed before shutting down nodes. This protocol can delay the

shutdown procedure when the timeouts do overlap just slightly or one node still works on the last OCDs candidates.

5 Conclusion

- Summary
- real conclusion
- future work
 - what’s not implemented yet (merging of results, deduplication, fixing algorithm)
 - optimize (sort and compare on multiple columns (consider data types or merge values to one or define ordering over indices?))
 - interesting questions: how to prove we have found all edgecases / formalize protocols
 - Test memory boundaries, possibly optimize memory consumption (delta protocol?), look at garbage collector overhead
 - Find and fix work stealing edge cases
 - evaluate timeouts
 - The actor model provides a dynamic aspect. This can be used to run multiple approaches in a single system and to exchange information between them. The goal of this would be to benefit from the advantages of the different approaches and at the same time decreasing their disadvantages. Our approach does not make use of this aspect of the actor system so far.

References

- [Co19] Consonni, C.; Sottovia, P.; Montresor, A.; Velegrakis, Y.: Discovering order dependencies through order compatibility. In: International Conference on Extending Database Technology. 2019.
- [GH83] Ginsburg, S.; Hull, R.: Order dependency in the relational model. Theoretical computer science 26/1-2, pp. 149–195, 1983.
- [Go19] Godfrey, P.; Golab, L.; Kargar, M.; Srivastava, D.; Szlichta, J.: Errata Note: Discovering Order Dependencies through Order Compatibility, 2019, arXiv: 1901.0205091 [cs.DB].
- [LN16] Langer, P.; Naumann, F.: Efficient order dependency detection. The VLDB Journal—The International Journal on Very Large Data Bases 25/2, pp. 223–241, 2016.

- [SGG12] Szlichta, J.; Godfrey, P.; Gryz, J.: Fundamentals of order dependencies. Proceedings of the VLDB Endowment 5/11, pp. 1220–1231, 2012.
- [Sz17] Szlichta, J.; Godfrey, P.; Golab, L.; Kargar, M.; Srivastava, D.: Effective and complete discovery of order dependencies via set-based axiomatization. Proceedings of the VLDB Endowment 10/7, pp. 721–732, 2017.