

A Lightweight Primitive for Online Tuning

by Tomer Kaftan (UW), Magdalena Balazinska (UW), Alvin Cheung (UW), Johannes Gehrke (Microsoft)



Logical Operators have multiple physical Operators... The system should automatically choose!

	A A A A A A A A A A A A A A A A A A A
1000	WIKIPEDIA The Free Encyclopedia

B no

Pages in category "Join algorithms"

The following 5 pages are in this category, out of 5 total.

в

Block nested loop

н

Hash join

Ν

Nested loop join

s

- Sort-merge join
- Symmetric Hash Join

Guideseri		alaga Sakaka Mg	a log a	e ²	logie, on svenge, ward case space complexity is 12. Redgework variation is logie, worst case.	place sort is not stable : stable weations axis.	Pations	Cubicater is usually done in-place with $O(\log r)$ stack space $\frac{\partial O(r)}{\partial r}$
	Verge acrit	alışın	a log a	alaga	s. A hybrid block marga son: Is D(1) mars.	Yes	Merging	Highly parameterize (up to Chog 4) using the Three Hangabane. Algorithm $^{[k]}$ or, mani practically, Coleia parallel merge acts for processing large encounts of data.
	toplace range sof	-	-	s log ² n Gee above, for hybriol frecis a loga	1	***	saya.	Currice implemented as a stable port based on stable in-place marging $\frac{14}{14}$
	Heapaon.	nlışa	alogs	a laga	1	Na	Selector	
	mention and	п	- 2 ²	n ²	1	716	Investion.	$\Omega(n+q)_{n}$ is the worst state over engineers that have a inversion.
	Intersort	alışın	nloga	alaga	log s.	No	Partitioning & Selection	Open in several 871, implementations.
	Selection spet.	e2	31	£ ²	1	Na	Selector.	Stable with $\Omega(t)$ extra space. For example using lists \overline{N}
	Treat	п	n log a	alaga	د	Yes	Intercion & Marging	Market Coordinations when the data is already coded or reverse source.
	Cubesol.	n	a loga	elage	٤	No	Inserters	Hakas it comparisons when the data is sineady sorted or revenue sorted.
	Shell sort	alaga	$s kg^2 s$ $s^{5/1}$	Depends on gap Asquinter; and Knownia. s.log ² o	ı	Na	Investory	Small sole state, so use of call state, meaning but, such when memory a strippen user such as encedence and other numeric approximate later case in the such work base may near sole actived togeths. With best case in the two transmission in M.
1	Dubble sort		*	1 ²	1	794	Exchanging	Tity code size.
	Minary tree eart	ninge	nikga	n logn (balanced)	τ	146	Interface.	When using a sall-balancing binary search thes.
1	Cycle soft	e ²	π^{T}	E,	1	Na	Investory	In place with theoretically optimal number of writes.
1	Library sort	- 10	a log s	6 ²		764	intertion	
	Padence sorting	п	-	winger		Na	Insertion & Selection	Finds all the longest increasing a data parameter in $O(t\log q)$
	Sinceficari.	п	n log a	a laga	3	No	Sector	An adaptive variant of hasport based upon the Leonards sequence rather than a traditional binary hasp.
	Stand and	п	- 2°	*		194	integral on	
1	Townshield set	ninge	u kiga	ninge	NIG	No	Sector	Variation of Heap Bol.
1	Codital sort	D.	4 ²	£ ²	1	79.0	Exchanging	
	Costo Apro	a been	°		1	09	Easteroing	Faster than table act on average.

State - Method

Other notes

The following 3 pages are in this category, out of 3 total. Thi

С

Cache-oblivious distribution sort

Е

F

- External sorting
- Funnelsort

Logical Operators have multiple physical Operators... The system should automatically choose!

(Some) Prior Work on Query Optimization

(Some) Prior Work on Query Optimization

- Static Query Optimizers
 - cardinality & selectivity estimation , heuristics, cost models

(Some) Prior Work on Query Optimization

- Static Query Optimizers
 - cardinality & selectivity estimation , heuristics, cost models
- Adaptive Query Optimization
 - Query re-optimization (update cardinality & selectivity estimates)
 - Adaptive operators (scans, aggregates, etc.)
 - Eddies & adaptive tuple routing (operator reordering)

• Designing good query optimizers takes time!

- Designing good query optimizers takes time!
- Requires deep knowledge of the operators and significant development effort

- Designing good query optimizers takes time!
- Requires deep knowledge of the operators and significant development effort
- Spark SQL took 2 years to go from heuristics-based optimization to cost-based optimization! ^[1]

- Designing good query optimizers takes time!
- Requires deep knowledge of the operators and significant development effort
- Spark SQL took 2 years to go from heuristics-based optimization to cost-based optimization! ^[1]
- Existing adaptive approaches just push the development overhead to physical execution

[1] http://databricks.com/blog/2017/08/31/cost-based-optimizer-in-apache-spark-2-2.html

- Designing good query optimizers takes time!
- Requires deep knowledge of the operators and significant development effort
- Spark SQL took 2 years to go from heuristics-based optimization to cost-based optimization!^[1]
- Existing adaptive approaches just push the development overhead to physical execution
- Modern data processing applications involve diverse, sophisticated operators, not just relational operators!

[1] http://databricks.com/blog/2017/08/31/cost-based-optimizer-in-apache-spark-2-2.html



"A Cuttlefish pretending to be a rock"

*Image Sourced from https://www.flickr.com/photos/silkebaron/32001215104



"A Cuttlefish pretending to be a rock"

*Image Sourced from https://www.flickr.com/photos/silkebaron/32001215104

Generate Training Data from:





*caption-generating model portion of the logical plan inspired by: Xu et al. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. ICML 2015 6



Diverse, sophisticated operators, with multiple physical alternatives!

*caption-generating model portion of the logical plan inspired by: Xu et al. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. ICML 2015 6

Example Operator: Convolution

Example Operator: Convolution

Tested 3 convolution algorithms on 8000 Flickr images



Can we optimize without a full-fledged optimizer?

Prior Work: Tuning Black-box Operators

Prior Work: Tuning Black-box Operators

- Offline Autotuning
 - Searches through arbitrary physical plan spaces
 - Requires representative workloads
 - Offline training time

Prior Work: Tuning Black-box Operators

- Offline Autotuning
 - Searches through arbitrary physical plan spaces
 - Requires representative workloads
 - Offline training time
- Micro Adaptivity in Vectorwise
 - Reinforcement learning chooses physical flavors of black-box vectorized operators
 - Limited to vectorized operators, does not explore multi-core settings

A Lightweight Primitive for Online Tuning



Workload developer (or the query optimizer) inserts calls to Cuttlefish's API to pick physical operators *during execution*

A Lightweight Primitive for Online Tuning



Workload developer (or the query optimizer) inserts calls to Cuttlefish's API to pick physical operators *during execution*

A Lightweight Primitive for Online Tuning



Workload developer (or the query optimizer) inserts calls to Cuttlefish's API to pick physical operators *during execution*

A Lightweight Primitive for Online Tuning



Developer maps tuning rounds to the execution model of each operator:

A Lightweight Primitive for Online Tuning



Developer maps tuning rounds to the execution model of each operator:

• Regex: One round per HTML Doc

A Lightweight Primitive for Online Tuning



Developer maps tuning rounds to the execution model of each operator:

- Regex: One round per HTML Doc
- Convolve: One round per image

A Lightweight Primitive for Online Tuning



Developer maps tuning rounds to the execution model of each operator:

- Regex: One round per HTML Doc
- Convolve: One round per image
- Parallel Distributed Join: One round per partition

- I. Problem & Motivation
- II. The Cuttlefish API
- III. Bandit-based Online Tuning
- IV. Distributed Tuning Approach
- V. Contextual Tuning
- VI. Handling Nonstationary Settings
- VII.Other Operators
- VIII.Conclusion

1. Construct a tuner (from a set of choices)

- 1. Construct a tuner (from a set of choices)
- 2. Tuner.choose (pick one of the choices)

- 1. Construct a tuner (from a set of choices)
- 2. Tuner.choose (pick one of the choices)
- 3. Tuner.observe (observe a reward for a choice)

- 1. Construct a tuner (from a set of choices)
- 2. Tuner.choose (pick one of the choices)
- 3. Tuner.observe (observe a reward for a choice)

Cuttlefish tuners maximize the total reward after multiple choose-observe tuning rounds
def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...

def loopConvolve(image, filters): ...

- def fftConvolve(image, filters): ...
- def mmConvolve(image, filters): ...

tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:

convolve, token = tuner.choose()

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:

convolve, token = tuner.choose()

start = now()
result = convolve(image, filters)
elapsedTime = now() - start

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])



convolve, token = tuner.choose()

start = now()
result = convolve(image, filters)
elapsedTime = now() - start
reward = computeReward(elapsedTime)

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])



convolve, token = tuner.choose()

start = now()
result = convolve(image, filters)
elapsedTime = now() - start
reward = computeReward(elapsedTime)

tuner.observe(token, reward)

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])



convolve, token = tuner.choose()

start = now()
result = convolve(image, filters)
elapsedTime = now() - start
reward = computeReward(elapsedTime)

tuner.observe(token, reward)

```
output result
```

Cuttlefish

- I. Problem & Motivation
- II. The Cuttlefish API

III. Bandit-based Online Tuning

- IV. Distributed Tuning Approach
- V. Contextual Tuning
- VI. Handling Nonstationary Settings

VII.Other Operators

VIII.Conclusion





Multi-armed Bandit Problem

Multi-armed Bandit Problem

• K possible choices (called arms)

Multi-armed Bandit Problem

- K possible choices (called arms)
- Arms have unknown reward distributions

Multi-armed Bandit Problem

- K possible choices (called arms)
- Arms have unknown reward distributions
- At each round: select an Arm and observe a reward

Multi-armed Bandit Problem

- K possible choices (called arms)
- Arms have unknown reward distributions
- At each round: select an Arm and observe a reward

Goal: Maximize Cumulative Reward (by balancing exploration & exploitation)



















• Gaussian runtimes with initially unknown means and variances



- Gaussian runtimes with initially unknown means and variances
- Belief distributions form t-distributions
 - Depend only on sample mean, variance, count



- Gaussian runtimes with initially unknown means and variances
- Belief distributions form t-distributions
 - Depend only on sample mean, variance, count
- No meta-parameters, yet works well for diverse operators



- Gaussian runtimes with initially unknown means and variances
- Belief distributions form t-distributions
 - Depend only on sample mean, variance, count
- No meta-parameters, yet works well for diverse operators
- Constant memory overhead, 0.03 ms per tuning round

• Prototype in Apache Spark

- Prototype in Apache Spark
- Tune between three convolution algorithms (Nested Loops, FFT, or Matrix Multiply)
 - Reward: -1*elapsedTime (maximizes throughput)

- Prototype in Apache Spark
- Tune between three convolution algorithms (Nested Loops, FFT, or Matrix Multiply)
 - Reward: -1*elapsedTime (maximizes throughput)
- Convolve 8000 Flickr images with sets of filters (~32gb)
 - Vary number & size of filters

- Prototype in Apache Spark
- Tune between three convolution algorithms (Nested Loops, FFT, or Matrix Multiply)
 - Reward: -1*elapsedTime (maximizes throughput)
- Convolve 8000 Flickr images with sets of filters (~32gb)
 - Vary number & size of filters
- Run on an 8-node (AWS EC2 4-core r3.xlarge) cluster.
 - 32 total cores, ~252 images per core

- Prototype in Apache Spark
- Tune between three convolution algorithms (Nested Loops, FFT, or Matrix Multiply)
 - Reward: -1*elapsedTime (maximizes throughput)
- Convolve 8000 Flickr images with sets of filters (~32gb)
 - Vary number & size of filters
- Run on an 8-node (AWS EC2 4-core r3.xlarge) cluster.
 - 32 total cores, ~252 images per core
- *Very* compute intensive
 - (Some configs up to 45 min on a single node)

Convolution Results



Relative throughput normalized against the highest-throughput algorithm

Convolution Results



Relative throughput normalized against the highest-throughput algorithm
Convolution Results



Relative throughput normalized against the highest-throughput algorithm

Cuttlefish

- I. Problem & Motivation
- II. The Cuttlefish API
- III. Bandit-based Online Tuning
- **IV. Distributed Tuning Approach**
- V. Contextual Tuning
- VI. Handling Nonstationary Settings
- VII.Other Operators
- VIII.Conclusion

- 1. Choosing and observing occur throughout a cluster
 - To maximize learning, need to communicate

- 1. Choosing and observing occur throughout a cluster
 - To maximize learning, need to communicate
- 2. Synchronization & communication overheads

- 1. Choosing and observing occur throughout a cluster
 - To maximize learning, need to communicate
- 2. Synchronization & communication overheads
- 3. Feedback delay
 - How many times is `choose' called before an earlier reward is observed?
 - Fortunately, theoretically sound to have delays

Centralized Tuner









Peer-to-Peer is also a possibility, but requires more communication







• When choosing: aggregate local & non-local state



- When choosing: aggregate local & non-local state
- When observing: update the local state



- When choosing: aggregate local & non-local state
- When observing: update the local state
- Model store aggregates non-local state

Results with Distributed Approach



Relative throughput normalized against the highest-throughput algorithm

Results with Distributed Approach



Throughput normalized against an ideal oracle that always picks the fastest algorithm

Results with Distributed Approach



Throughput normalized against an ideal oracle that always picks the fastest algorithm

Cuttlefish

- I. Problem & Motivation
- II. The Cuttlefish API
- III. Bandit-based Online Tuning
- IV. Distributed Tuning Approach
- V. Contextual Tuning (by learning cost models)
- VI. Handling Nonstationary Settings
- VII.Other Operators
- VIII.Conclusion





- Best physical operator for each round may depend on current context
 - e.g. convolution performance depends on the image & filter dimensions



- Best physical operator for each round may depend on current context
 - e.g. convolution performance depends on the image & filter dimensions

- Users may know important context features
 - e.g. from the asymptotic algorithmic complexity



- Best physical operator for each round may depend on current context
 - e.g. convolution performance depends on the image & filter dimensions

- Users may know important context features
 - e.g. from the asymptotic algorithmic complexity
- Users can specify context in Tuner.choose

 Linear contextual Thompson sampling learns a linear model that maps features to rewards

- Linear contextual Thompson sampling learns a linear model that maps features to rewards
- Feature Normalization & Regularization
 - Increased robustness towards feature choices

- Linear contextual Thompson sampling learns a linear model that maps features to rewards
- Feature Normalization & Regularization
 - Increased robustness towards feature choices
- Effectively learns a cost model

Tuning Convolution with Cuttlefish

def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...

tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:

convolve, token = tuner.choose()

```
start = now()
result = convolve(image, filters)
elapsedTime = now() - start
reward = computeReward(elapsedTime)
```

tuner.observe(token, reward)

```
output result
```

Tuning Convolution with Cuttlefish



Contextual Convolution Results



Throughput normalized against an ideal oracle that always picks the fastest algorithm

Cuttlefish

- I. Problem & Motivation
- II. The Cuttlefish API
- III. Bandit-based Online Tuning
- IV. Distributed Tuning Approach
- V. Contextual Tuning

VI. Handling Nonstationary Settings

VII.Other Operators

VIII.Conclusion

- Runtimes may drift over time, or differ across nodes
 - heterogenous cluster, changing resource availabilities, data properties varying throughout the workload, etc.
 - E.g. web crawl data and images may be stored sorted by website. This could correlate with performance

- Runtimes may drift over time, or differ across nodes
 - heterogenous cluster, changing resource availabilities, data properties varying throughout the workload, etc.
 - E.g. web crawl data and images may be stored sorted by website. This could correlate with performance
- We might not be capturing sufficient context!

- Runtimes may drift over time, or differ across nodes
 - heterogenous cluster, changing resource availabilities, data properties varying throughout the workload, etc.
 - E.g. web crawl data and images may be stored sorted by website. This could correlate with performance
- We might not be capturing sufficient context!
- Standard multi-armed bandit techniques fail
- Prior work: dynamic bandit approaches
 - Sliding windows, discounting older observations, reset on change detection, etc.
 - Good for dealing with changes over time

- Prior work: dynamic bandit approaches
 - Sliding windows, discounting older observations, reset on change detection, etc.
 - Good for dealing with changes over time
- Prior work: bandit clustering approaches
 - identify & share learning among agents solving similar bandit problems
 - Good for dealing with differences between cores

- Prior work: dynamic bandit approaches
 - Sliding windows, discounting older observations, reset on change detection, etc.
 - Good for dealing with changes over time
- Prior work: bandit clustering approaches
 - identify & share learning among agents solving similar bandit problems
 - Good for dealing with differences between cores
- Need dynamic bandit clustering where agents' underlying problems may change over time!





Observations



41



To Lower Overheads



To Lower Overheads



To Lower Overheads

Observations

Store only one 'aggregated old state' per epoch At epoch end: If similar to old, merge into 'old state' . Otherwise, replace 'old state'

Identify (& merge) similar non-local states only at communication rounds, in the centralized model store

Nonstationary Results



Throughput normalized against an ideal oracle that always picks the fastest algorithm

Cuttlefish

- I. Problem & Motivation
- II. The Cuttlefish API
- III. Bandit-based Online Tuning
- IV. Distributed Tuning Approach
- V. Contextual Tuning
- VI. Handling Nonstationary Settings
- **VII.Other Operators**
- VIII.Conclusion





- Tune between four regular expression searching libraries
 - Built-in Java Regex and 3 third-party libraries



- Tune between four regular expression searching libraries
 - Built-in Java Regex and 3 third-party libraries
- Search through 256k Common Crawl docs (~30gb uncompressed)
 - one tuning round per doc

Regex Operator

- Tune between four regular expression searching libraries
 - Built-in Java Regex and 3 third-party libraries
- Search through 256k Common Crawl docs (~30gb uncompressed)
 - one tuning round per doc
- Test 8 Regexes sourced from regex-sharing website RegExr
 - Match hyperlinks, trigrams, valid emails, color codes, etc.

Regex Operator

- Tune between four regular expression searching libraries
 - Built-in Java Regex and 3 third-party libraries
- Search through 256k Common Crawl docs (~30gb uncompressed)
 - one tuning round per doc
- Test 8 Regexes sourced from regex-sharing website RegExr
 - Match hyperlinks, trigrams, valid emails, color codes, etc.
- Multiple of orders of magnitude variation in performance
 - Email validation regex w/ built-in java utilities takes 33µs to process the fastest document, but over 1000s for the slowest document

Regex Operator

- Tune between four regular expression searching libraries
 - Built-in Java Regex and 3 third-party libraries
- Search through 256k Common Crawl docs (~30gb uncompressed)
 - one tuning round per doc
- Test 8 Regexes sourced from regex-sharing website RegExr
 - Match hyperlinks, trigrams, valid emails, color codes, etc.
- Multiple of orders of magnitude variation in performance
 - Email validation regex w/ built-in java utilities takes 33µs to process the fastest document, but over 1000s for the slowest document
- 8-node (AWS EC2 4-core r3.xlarge) cluster

Regex Results



Note: Y-axis is Log-scale

• Hash-partition relations according to join attributes

- Hash-partition relations according to join attributes
- On each partition, pick a local hash join or a local sort-merge join

- Hash-partition relations according to join attributes
- On each partition, pick a local hash join or a local sort-merge join
- Rewards capture total join time
 - measure from when joins begin until result iterators are fully consumed

- Hash-partition relations according to join attributes
- On each partition, pick a local hash join or a local sort-merge join
- Rewards capture total join time
 - measure from when joins begin until result iterators are fully consumed
- Set as Spark SQL 2.2's join for all equijoins too large to broadcast
 - No heuristics and cost models in the query optimizer, falls back on explicit configurations (defaults to global sort-merge join)

- Hash-partition relations according to join attributes
- On each partition, pick a local hash join or a local sort-merge join
- Rewards capture total join time
 - measure from when joins begin until result iterators are fully consumed
- Set as Spark SQL 2.2's join for all equijoins too large to broadcast
 - No heuristics and cost models in the query optimizer, falls back on explicit configurations (defaults to global sort-merge join)
- Test on TPC-DS benchmark (scale factor 200)

- Hash-partition relations according to join attributes
- On each partition, pick a local hash join or a local sort-merge join
- Rewards capture total join time
 - measure from when joins begin until result iterators are fully consumed
- Set as Spark SQL 2.2's join for all equijoins too large to broadcast
 - No heuristics and cost models in the query optimizer, falls back on explicit configurations (defaults to global sort-merge join)
- Test on TPC-DS benchmark (scale factor 200)
- Configure queries to use 512 shuffle / join partitions

Join Results (Query Throughput)



Join Results (Query Throughput)



But, requires exploration & provides no 'special ordering' benefits

Join Results (Query Throughput)



Cuttlefish join usually faster (Join throughput graphs even more dramatic) But, requires exploration & provides no 'special ordering' benefits

Cuttlefish

- I. Problem & Motivation
- II. The Cuttlefish API
- III. Bandit-based Online Tuning
- IV. Distributed Tuning Approach
- V. Contextual Tuning
- VI. Handling Nonstationary Settings
- VII.Other Operators

VIII.Conclusion

Cuttlefish

- A simple, flexible API for online tuning
- Thompson-sampling based tuning algorithms
- Supports contextual tuning (learns cost models)
- Distributed learning between workers
- Adapts to nonstationary workloads
- Prototyped in Apache Spark & successfully tunes convolution, regex, and join operators