

XMLディスクとその応用事例

高度メディア蓄積・管理手法研究グループ (奈良先端科学技術大学院大学 筑波大学 同志社大学)

研究の背景

近年のXMLデータの急速な普及

バイオデータベース、Webサービス、Wikiペディア、...

ホスト計算機でのXMLデータ処理および管理コスト

ストレージ側にXMLデータ処理をオフロード

研究目的

XMLディスクの開発

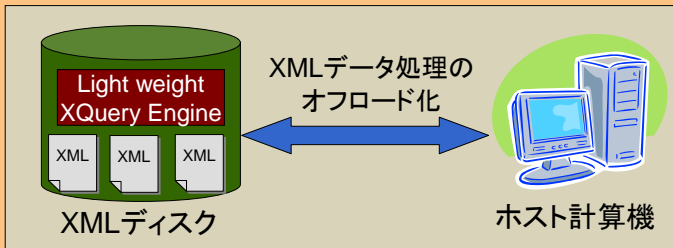
XQueryエンジンの軽量化と自律ディスクへの組み込み

XMLディスクの応用例

バイオ情報データの検索

XMLディスク

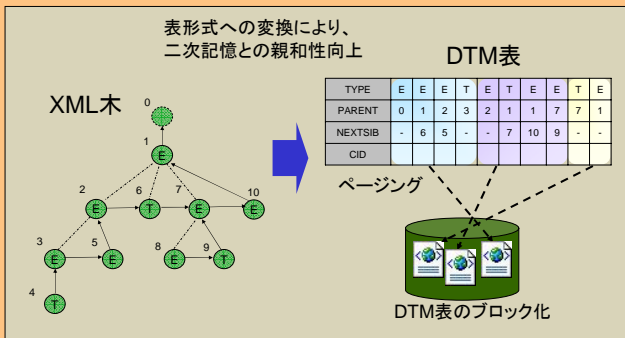
コンセプト



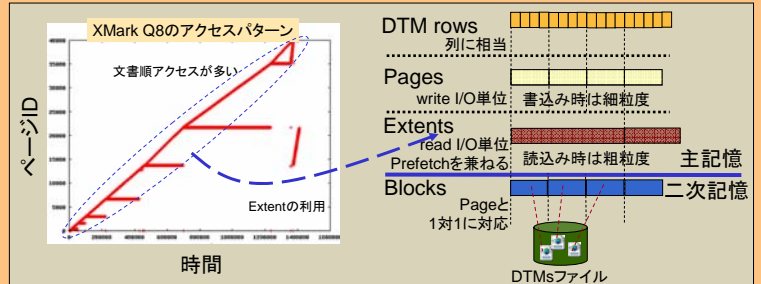
XQueryエンジンの省メモリ化

従来のXQueryエンジンは大量のメモリを消費

Document Table Model (DTM)を二次記憶への格納に拡張



I/O粒度の最適化 XQueryのページアクセスパターンからI/O粒度の決定



性能評価 XMarkベンチマークによる評価

Q	XMark13MB (ScaleFactor: 1)			XMark340MB (ScaleFactor:3)		
	DTM	可変長DTMs	固定長DTMs	DTM	可変長DTMs	固定長DTMs
Q1	9.838	3.531	3.609	9.875	28.84	10.704
Q2	9.822	4.296	3.891	10.203	28.593	11.75
Q3	9.89	5.219	4.887	10.328	30.281	14.359
Q4	9.812	5.453	5.288	10.375	29.841	15.282
Q5	9.203	3.719	3.983	10.25	29.172	10.14
Q6	10.988	10.156	8.578	9.908	32.825	28.375
Q7	12.989	15.582	14.42	11.047	39.203	61.908
Q8	10.872	8.857	5.515	11.344	30.5	14.875
Q9	10.853	8.837	9.25	914.872	31.31	25
Q10	18.813	23.157	27.422	13.422	49.704	78.453
Q11	18.328	18.907	18.938	786.344	99.887	87.287
Q12	18.703	12.984	12.658	351.812	67.583	80.158
Q13	9.781	3.437	3.563	10.719	28.781	9.84
Q14	11.719	12.583	11.484	10.781	34.47	37.822
Q15	9.785	3.25	3.158	9.828	28.989	8.578
Q16	9.812	3.38	3.381	10.5	28.15	9.408
Q17	10.109	4.158	3.884	10.468	28.891	11.83
Q18	10.25	4.989	4.75	10.872	28.822	13.422
Q19	11.532	9.594	9.787	11.5	33.857	28.922
Q20	10.312	5.344	6.158	10.313	30.219	14.785

heapサイズの影響 (ScaleFactor:10)

Q	heapサイズの影響 (ScaleFactor:10)	
	1024m	256m
Q1	37.82	38.858
Q2	40.884	41.825
Q3	38.858	39.31
Q4	41.841	43.888
Q5	30.489	31.718
Q6	93.32	88.187
Q7	151.704	153.738
Q8	50.25	44.578
Q9	101.882	110.884
Q10	328.531	400.875
Q11	731.344	747.888
Q12	488.781	501.63
Q13	33.887	33.547
Q14	152.437	143.25
Q15	32.234	27.32
Q16	27	27.36
Q17	31.14	32.47
Q18	38.235	38.14
Q19	138.47	ERR
Q20	58.219	51.453

大多数のクエリにおいてSaxon-SAより高速
Q9,Q11,Q12ではより顕著

メモリが少なくても安定して動作
(Saxon-SAは双方ともメモリ不足エラー)

デモンストレーション

複数のバイオデータベースの情報統合

ウェット実験者のための検索サポート
アノータのデータベース参照支援

バイオデータベースの問題点

同一塩基、タンパク質でもデータベースで異なるIDが使用
複数のデータベースとの情報統合が困難

デモ内容

BLASTで得られたアミノ酸配列の相同性から
複数のデータベースの柔軟な情報統合を行う

