

HTTP-FUSE KNOPPIX

須崎有康，八木豊志樹，飯島賢吾¹⁾，丹英之²⁾
産業技術総合研究所¹⁾，アルファシステムズ²⁾

はじめに

- HTTP-FUSE KNOPPIX とは
- ”KNOPPIX” の仕組みのおさらい
- 分割圧縮ファイルによるループバックデバイス
- Internet 対応方法
- 関連研究、今後の予定

HTTP-FUSE KNOPPIX とは

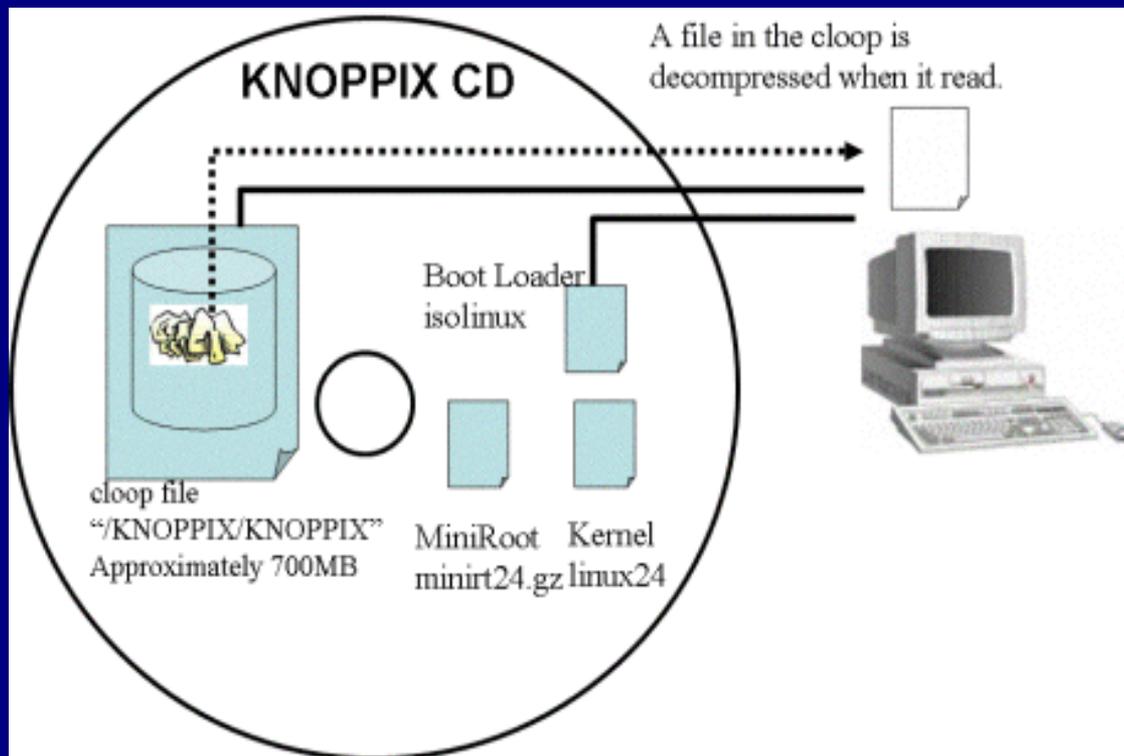
- 「Internet から起動する KNOPPIX」
- 大目的（予算要求用看板）
 - 現在の OS インストール / ブート方法の改革
 - OS インストールでは同じもののコピー
 - 内蔵デバイスに OS を固定するのは止めよう
 - Internet のどこかに OS はある
- 小目的（実利）
 - KNOPPIX のカスタマイズで CD を作り直すのは止めよう。

KNOPPIXのおさらい

- ブート時のデバイスの自動認識 / ドライバの組み込み機能 (Autoconfig) に優れている
 - anonymous user がネットワークブートするための必要条件
- ルートファイルシステムが圧縮ループバックデバイス (cloop) に格納され、扱いやすい。
 - ただし 1 巨大ファイルで、更新が面倒臭い。

KNOPPIX の特徴 (cloop)

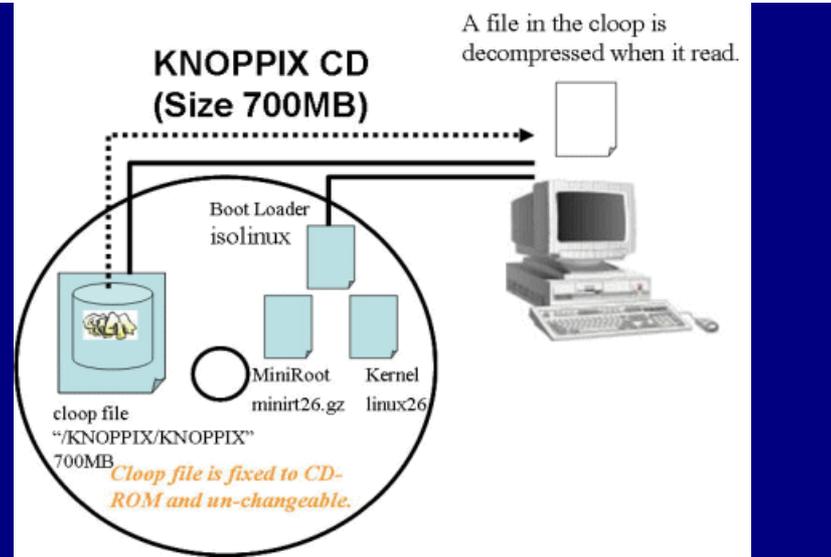
- 700M CD-ROM に圧縮ループバックデバイス CLOOP を使い、2Gに拡張。
- 圧縮にはzlibを利用。



cloop の改良

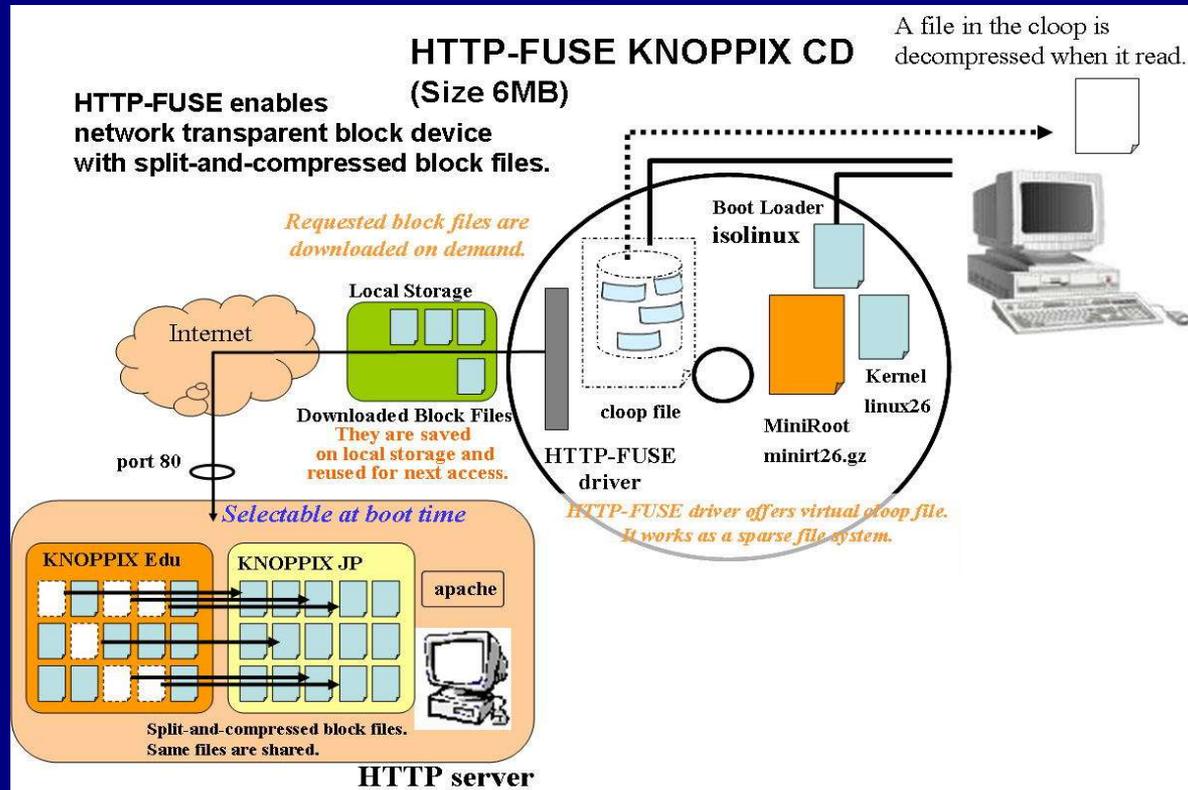
- cloop 問題点
 - 1 巨大ファイル (700MB) のループバックマウント
 - 700MB ダウンロードは辛い。
 - 一部の変更でも 700MB の作りなし。
- Internet から直接 cloop を利用できないか
 - 利用した部分だけをダウンロードし保存。
 - リモートとローカルの透過的 cloop
 - 差分更新を可能にする。
 - 細かいブロックファイルのループバックマウント
 - 脱 Client/Server。 P2P 的にばらまく。

- 現状の **KNOPPIX**
- 700MB CD にすべてを含む
 - boot loader, kernel, miniroot, RootFileSystem



- **HTTP-FUSE**
KNOPPIX

- 6MB CD
 - boot loader, kernel, miniroot
- rootFS は Internet より取得する
 - **Selectable!**

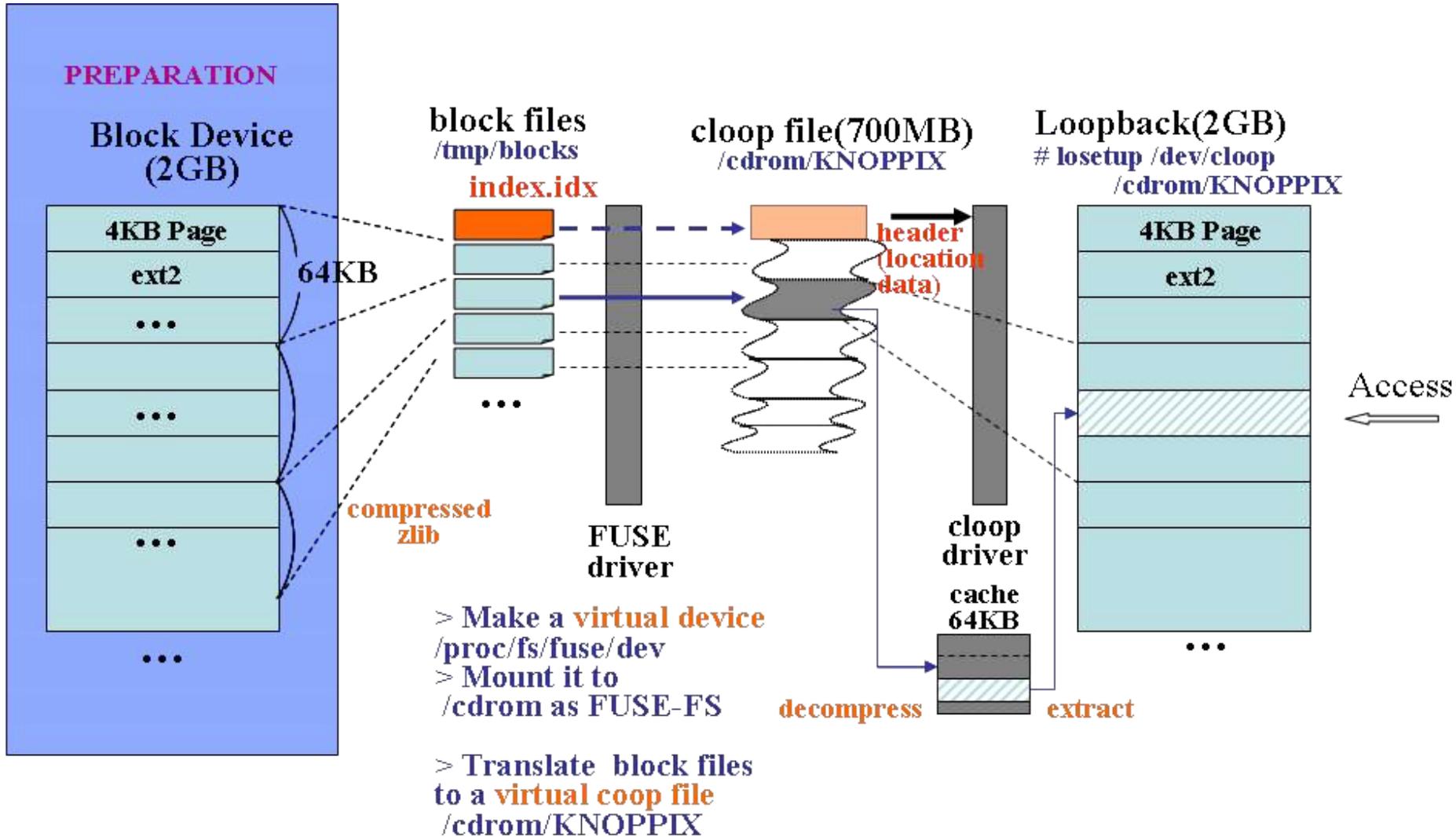


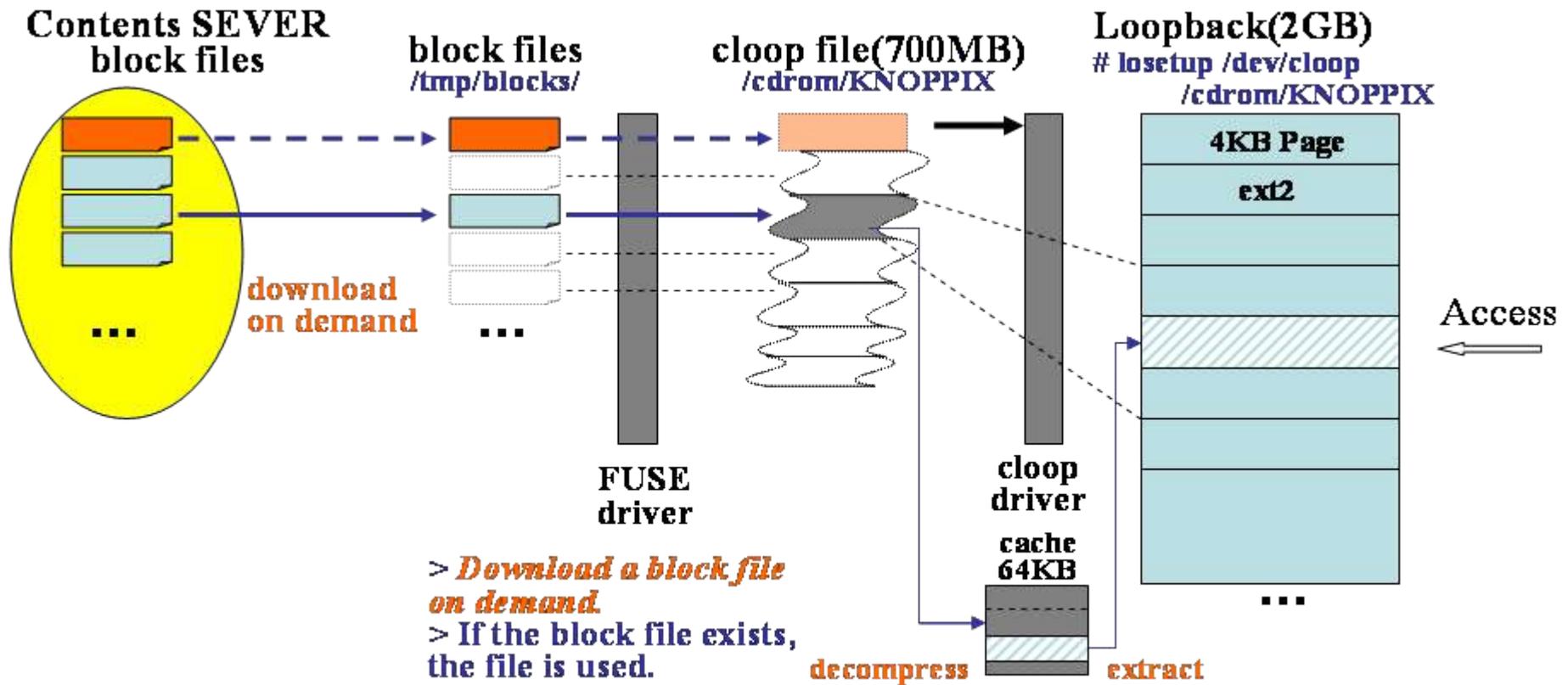
圧縮分割ファイルによる ループバックデバイス

- cloop を改良してリモートとローカルで透過的に扱えるようにする。
 - cloop は 64KB 毎のブロックを切り出し、圧縮して一つのループバックデバイスファイルを作っていた
 - 大きなまま扱わなければならないことが問題
 - 64KB 毎のブロックを切り出し、圧縮データを小さな**ブロックファイル**として扱い、それらをまとめてループバックデバイスとして扱えるようにする
 - ブロックファイルは**必要に応じて download** する。

圧縮分割ファイルによる ループバックデバイス

- ファイル名は md5sum の値。
 - 内容確認ができる。
 - 同一内容は 1 ファイル。
- ブロックファイルをまとめる手段として仮想ファイルシステム FUSE を利用
 - 分割圧縮ファイルから構成できよう wrapper を作成
 - FUSE を通してブロックファイル群を cloop ファイルとして構成する。





- > *Download a block file on demand.*
- > If the block file exists, the file is used.
- > *block files are erasable.*
- > Volume of /tmp/blocks/ keeps limit.

decompress extract

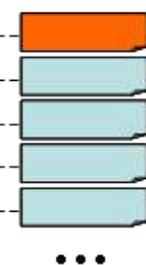
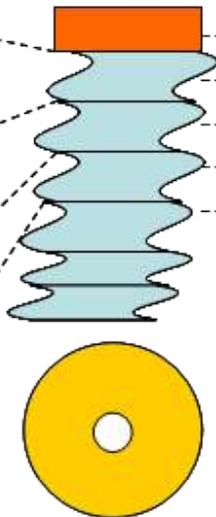
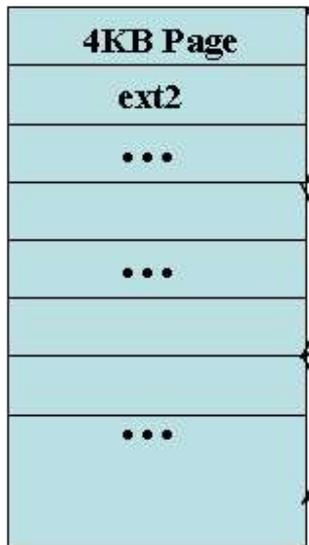
Block Device (2GB)

CD style

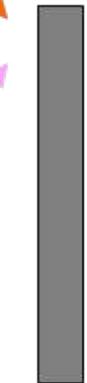
block file style

cloop file

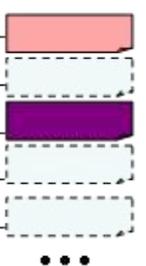
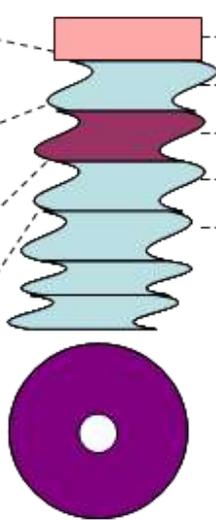
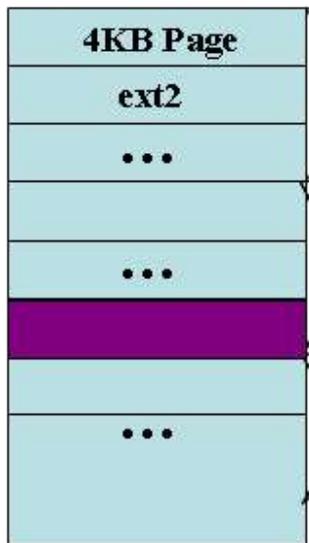
block files named by MD5



Same files
Reusable for FUSE



Update
apt-get install ...



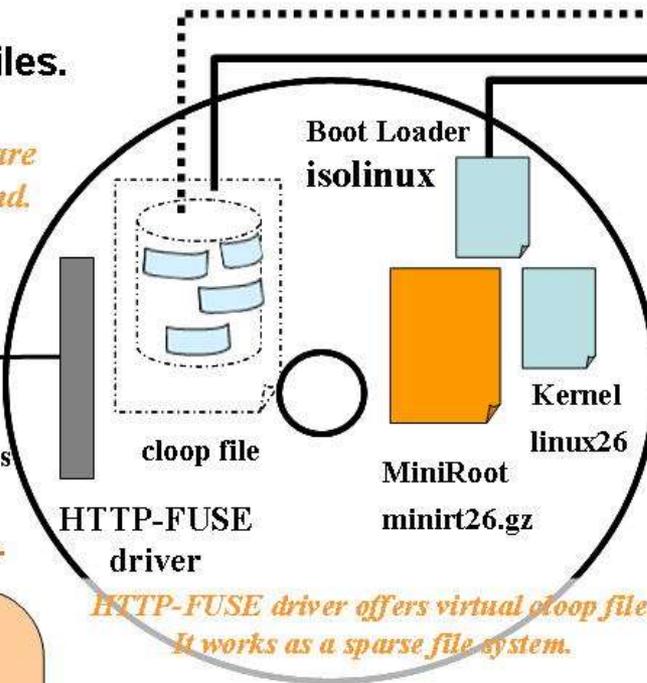
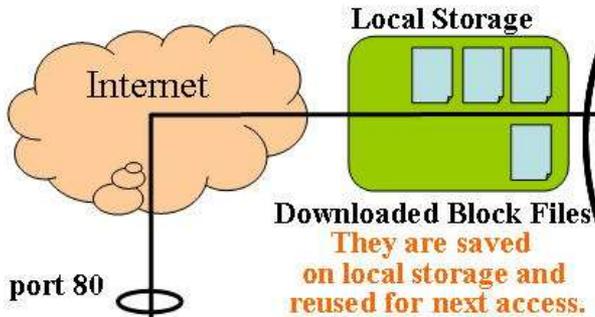
New KNOPPIX

HTTP-FUSE KNOPPIX CD (Size 6MB)

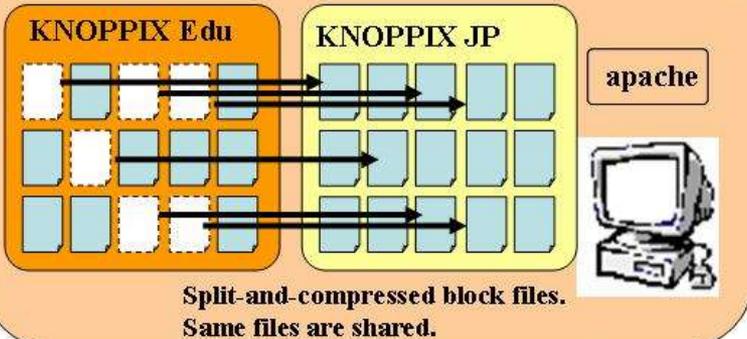
A file in the cloop is decompressed when it read.

HTTP-FUSE enables network transparent block device with split-and-compressed block files.

Requested block files are downloaded on demand.



Selectable at boot time



HTTP server

FUSE+CLOOP の特徴

- Pros

- リモート／ローカルの透過的に扱える
 - すべてダウンロードする必要なし
 - ローカルデバイスにダウンロードしたブロックファイルは再利用可能
- 更新があった場合、変更があったブロックに対応する分割ファイルとヘッダ情報 (index.idx) のみでよい。

- Cons

- 性能は？

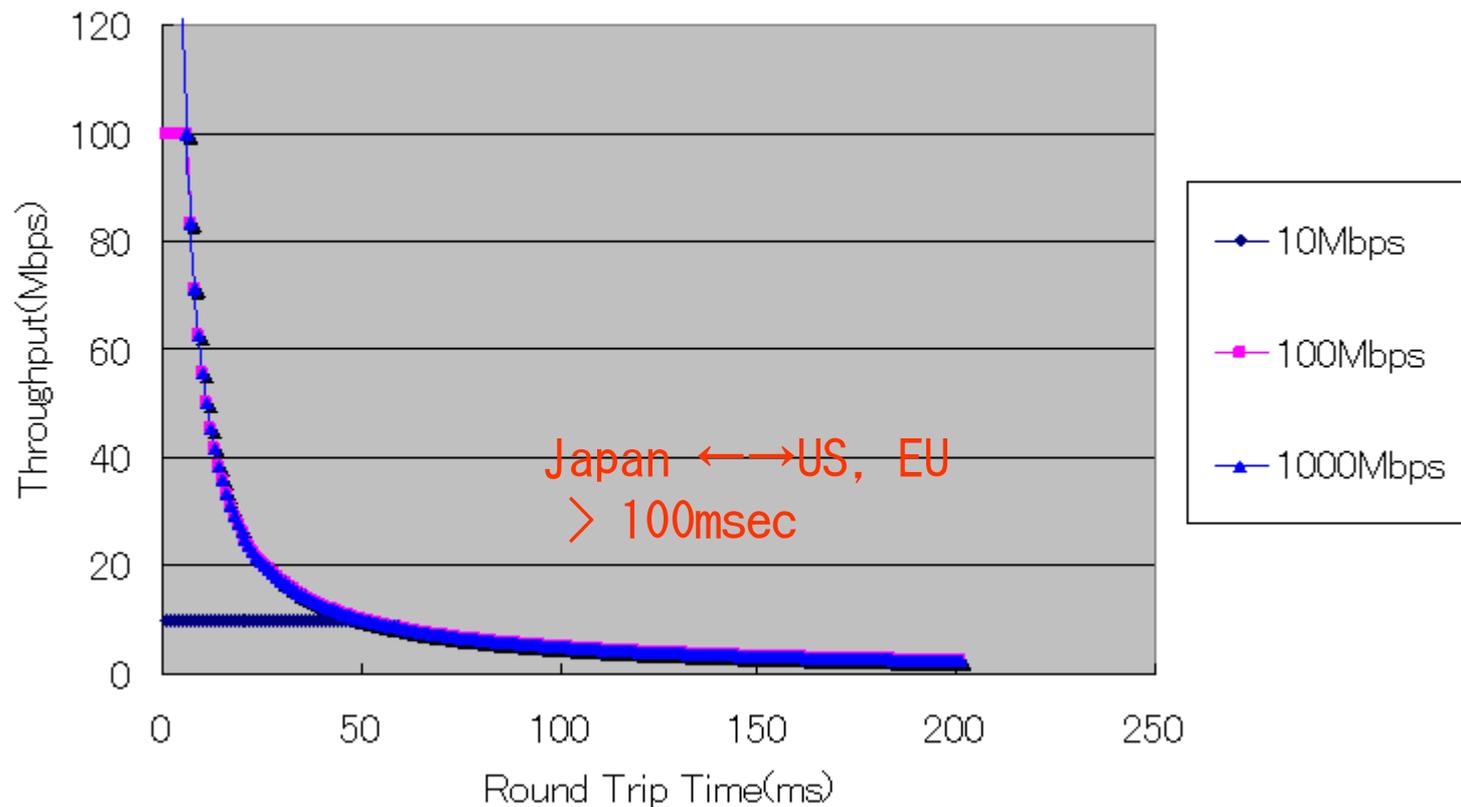
分割圧縮ブロックファイルの特徴

- 条件
 - (64KB cut) ブロックファイル
 - 31,206 files(Max:65,562 Min:84 Ave:21,801)
 - ブート時にダウンロードされるファイル
 - 4,424 files(87MB)
 - ダウンロードパターン(ext2へのアクセスパターン)
 - on demand (非連続)
- 要求
 - 小さいファイルを素早く、on demand 要求

どのようにしてブロックファイルを要求するか

- ダウンロード要求手順
 - ext2 → cloop → fuse → (Ineternet) → file
- cloop で access request が逐次化
 - コネクションが1つ。
 - NFS で使われているような多重 access request によるバンド幅拡張ができない。
 - レンテンシの影響を受けやすい。
 - Window size(64KB) 以下より小さいファイルでは性能がでない。
 - 64KB 切出しでは (Max:65,562 Min:84 Ave:21,801)

Latency and Throughput



Client の TCP windows size を 64KB, 1 connection とした場合。

10Mbps : $x < 50$ $y=10$, $x \geq 50$, $y = (50/x) * 10$

100Mbps: $x < 5$ $y=100$, $x \geq 5$, $y = (5/x) * 100$

1000Mbps: $x < 0.5$ $y=1000$, $x \geq 0.5$, $y = (0.5/x) * 1000$

切り出しブロックサイズの変更

- FUSE+CL00P では切り出しブロックサイズの変更は自由にできる。

	Files	Max	Min	Ave	down(sizeMB)	Full
64KB	31,206	65,562	84	21,801	4424(87)	83%
128KB	15,611	131,118	149	43,041	2876(115)	72%
256KB	7,821	262,230	277	85,372	2089(171)	59%
512KB	3,927	524,454	531	169,501	1436(231)	45%

同一内容のファイルが存在する場合はファイル数が減る

full はブートに実際使われる ext2 ページ割合

切り出しブロックサイズの変更

- FUSE での読み出し時間（オーバーヘッド）の影響
 - `dd of=/cdrom/KNOPPIX if=/dev/null (680MB)`
- bootchart 時間（init の終から KDE の起動前まで）
 - HTTP down load で起動。参考：CD 起動の場合 2:15 程度

	RAM-Disk	HTTP	bootchart 時間 (データ/ファイル数)
64KB	9.96	1:58	0:43(87MB,4224)
128KB	10.59	1:42	0:57(115MB,2876)
256KB	20.91	1:43	1:01(171MB,2089)
512KB	54.34	2:12	1:16(231MB,1436)

#ThinkPADT42, PentiumM 1.8G,100MNIC,HTTP は直結 RTT< 0.2msec
RTT の長い場合の影響は調べる必要あり。

ダウンロードの方法

- 高速にファイル（コピー）をばらまきたい。
- P2P? 特に BitTorrent
 - 通常の P2P はファイル検索機能がある。該当ファイルを探すのには適するが非常に遅い。
 - 残念ながら BitTorrent は大容量ファイルに適した設計。
 - ファイルをピースに分割し、複数のノードからピースを適当な順番でダウンロード（バンド幅拡張）し、最後に再構成することでバンド幅を稼いでいる。
- on demand に小さいファイルを多くダウンロードするのに適していない。

ダウンロードの方法

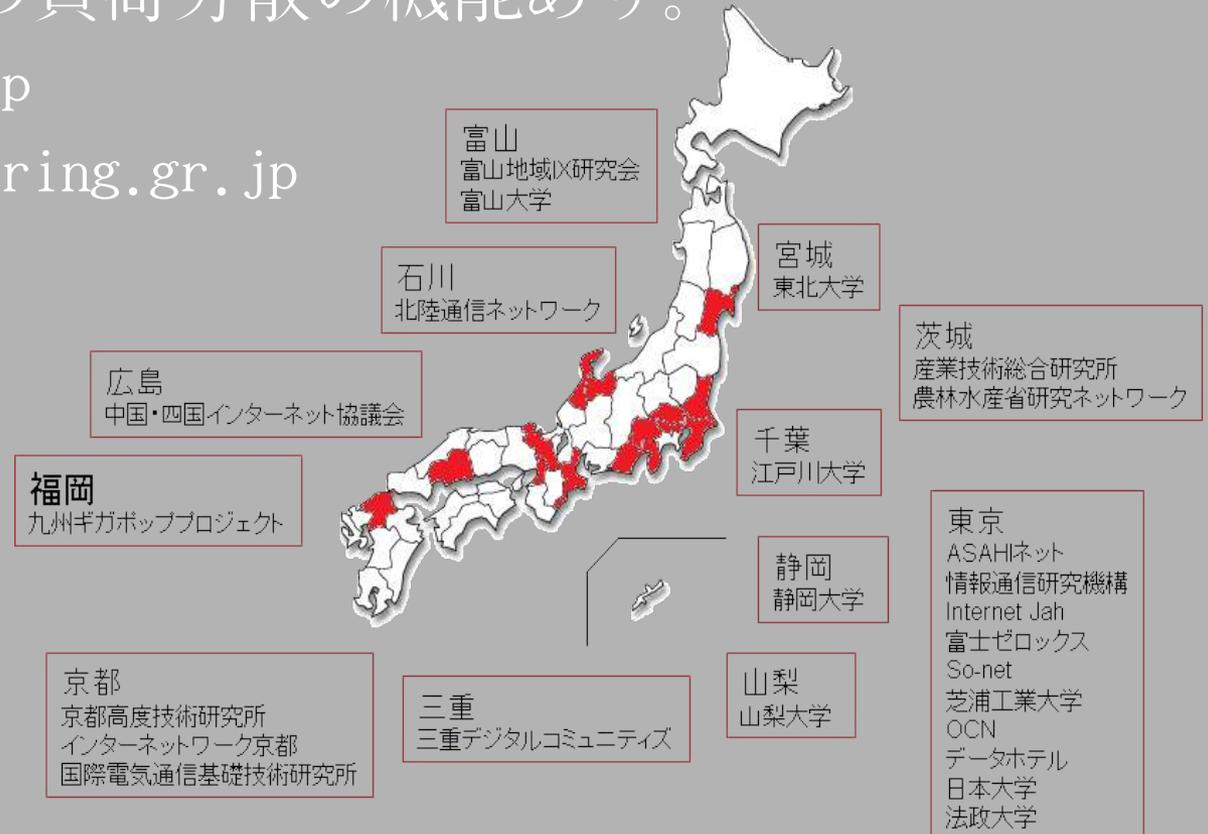
- CDN(Content Delivery Network)の利用。
 - 1つのURLで複数の配信法を有するもの。
 - ルーティング技術、キャッシュ、コンテンツ管理を行ない、効率的配信を実現する(仮想)ネットワーク。
 - 技術的蓄積や利用可能な資源がある。
 - RTT, レスポンスを短くする技術
 - BGP を利用したもの。サーバの負荷分散。
 - ミラーサーバ
 - proxy
 - P2P proxy(coral, P2P squid, dijjer,...)
 - reverse proxy
 - プロトコルは手軽で実績のある http

CDN の候補

- ミラーサーバ群の ring プロジェクト
 - 国内に 20 以上配置。
- P2P proxy の coral プロジェクト
 - PlanetLab を使い、世界中に配置。

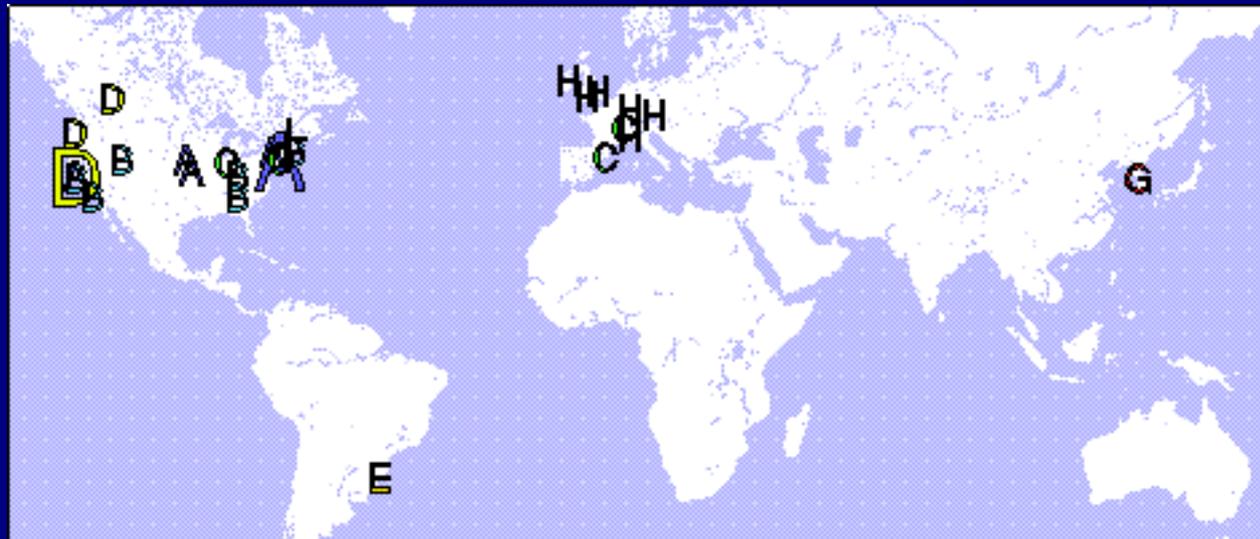
ring プロジェクト

- 日本のフリーソフト配信の最大ミラーサーバ
- tenbin, DNS balance によるクライアントへ近いミラー選択、サーバの負荷分散の機能あり。
 - www.t.ring.gr.jp
 - www.dnsbalance.ring.gr.jp



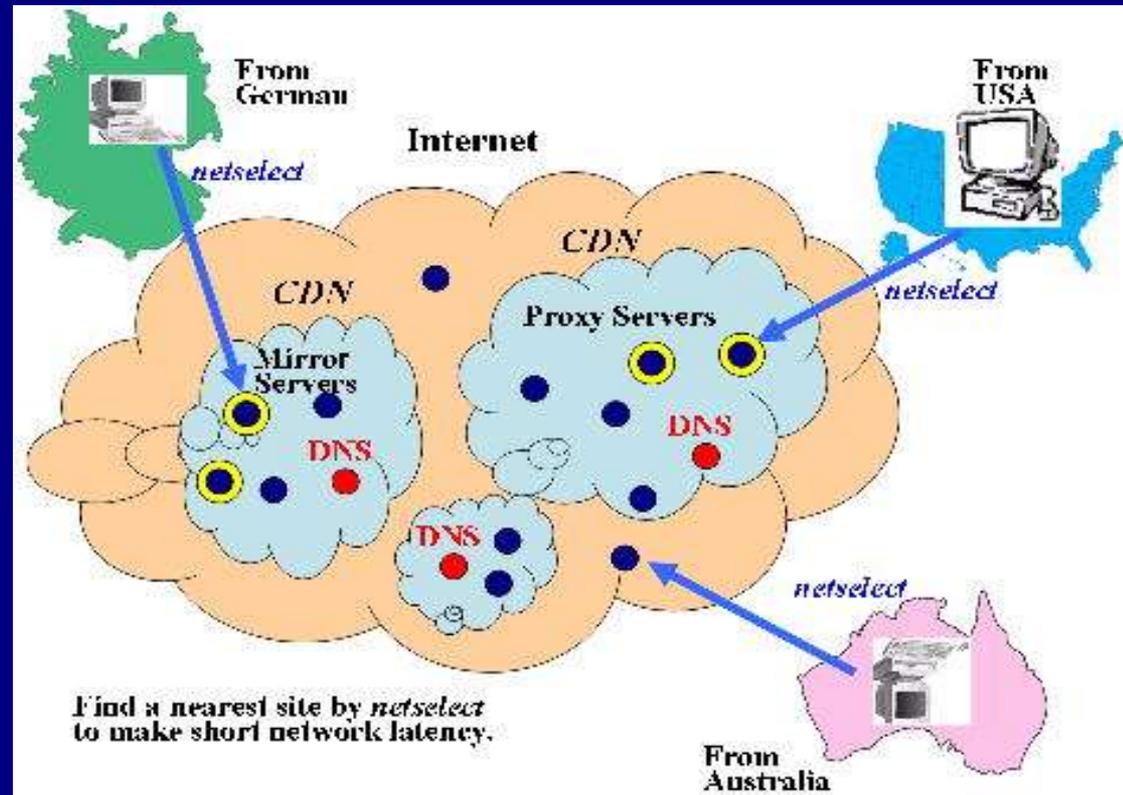
coral (P2P proxy)

- .nyud.net:8090 を URL に付けるのみ。
 - Ex: `www.aist.go.jp.nyud.net:8090/index.html`
 - nyud.net で名前解決
- DSHT(Distributed Sloppy Hash Table) によりホットスポットを回避
- coral 内にファイルがキャッシュされていればキャッシュを利用。



CDN とのネゴシエーション

- CDN のサーバに頼るだけでなく、クライアント側でも RTT が短くなるような CDN を選択。
- CDN とクライアントの相互のネゴシエーションによりダウンロードが決定。
 - 具体的には CDN が返す IP アドレス群に netselect をかけて選択。



現状

- 以上述べた機能をいれた HTTP-FUSE **KNOPPIX** を (04/14) 公開。
 - ブート環境 (6MB CD-ROM, USB ブート)
 - ブートオプション
 - http_proxy=
 - proxy を指定可能
 - memcache
 - 二記憶に保存せず、RAM DISK のみを使う

関連研究

- Plan9 の Venti [USENIX02]
 - アドレスではなく、書き込むブロック内容のハッシュで管理するストレージシステム
 - 内容が同じブロックは同一
 - ファイルシステムは別 (fossil)
- Locality awareness がある分散ファイルシステム
 - AFS, Coda, Shark [NSDI05] など。ファイルシステムレベルで管理。メタデータのキャッシュなど機能分割が明確
 - 書き込み可能。ネットワークの離脱、再接続で一貫性管理あり。
 - Cons: サーバ/クライアント。特殊なプロトコル。

関連研究

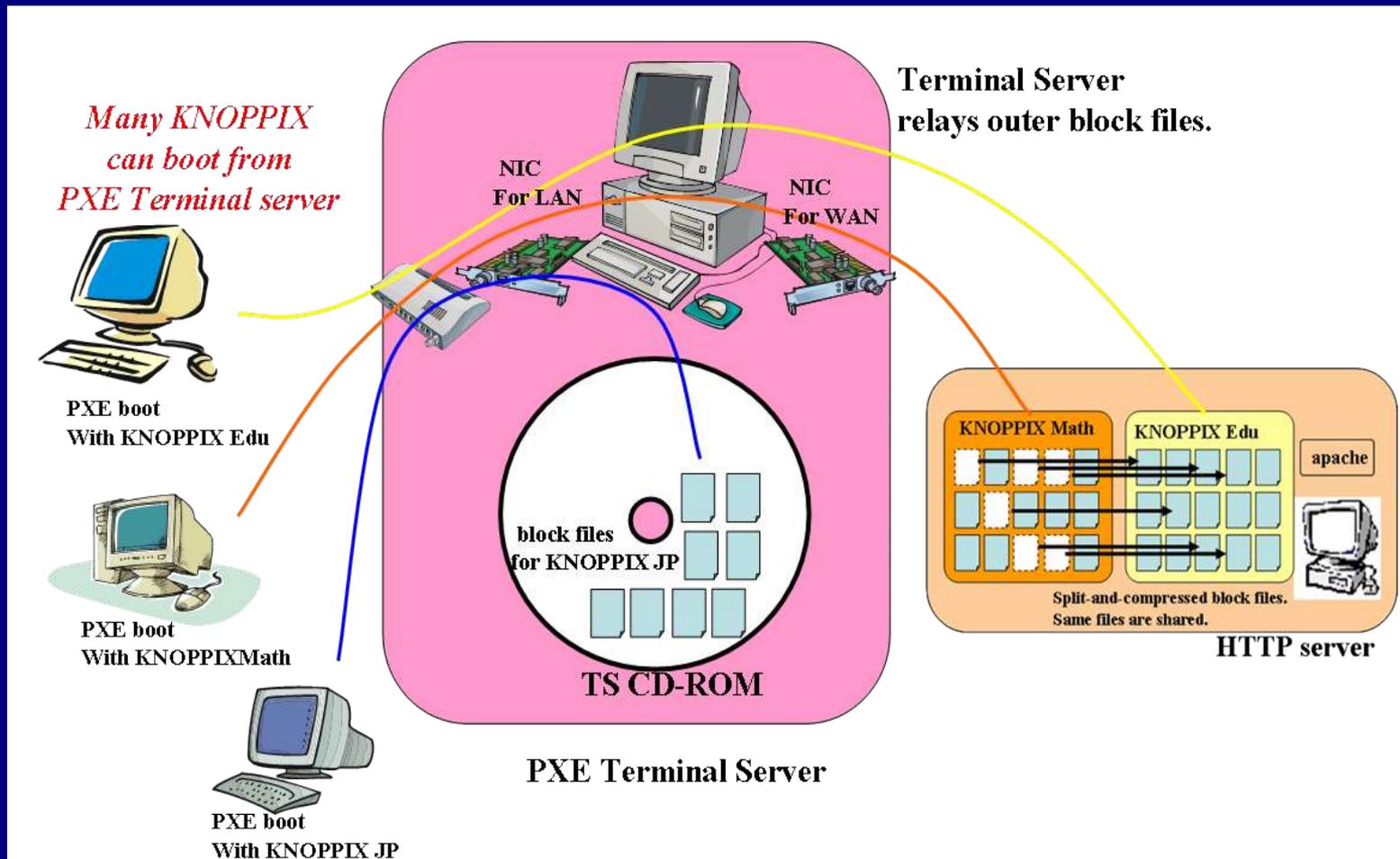
- UNIONFS
 - KNOPPIX3.8から導入された更新可能ファイルシステム
 - CD 上で apt-get が利用可能。
- CopyOnWrite を用いたオーバーレイブロックデバイス (丹、Linux Conference'05)
 - ブロックレベルの差分更新記録システム
 - 分割ブロックファイルに変更できれば、カスタマイズが容易になる

今後の課題／展望

- ブートデバイス問題
 - 現状ではBIOSのブートデバイス(CD,HD)を使う必要あり
 - Intelが中心に策定しているBIOS規格EFI(Extensible Firmware Interface)に期待
 - OSローダが自由に組み込めるようになり、今までの固定ブートデバイスから解放
- ファイルシステムの最適化
 - ブート時に使われるブロックを連続配置はできないか
 - フラグメンテーションの抑制。
- パッケージインストールにも対応したい。
 - 差分更新できないか。(関連: klik or zeroinstall)

PXE Terminal Server for HTTP-FUSE KNOPPIX

- N+I(6/8-10) で展示



まとめ

- 分割圧縮ブロックファイルによるループバックデバイスの作成。
- ネットワーク等価的にブロックファイルの利用が可能。
- ブロックファイルはコピー配布で構わない。現状ではCDNを使って配布。
 - この部分は改良の余地あり。
- この上で起動する **KNOPPIX** を作成。
 - HTTP-FUSE **KNOPPIX**