

LTU technologies



Image mining technologies and industrial challenges

Sébastien GILLES, Ph.D.
Chief scientist & co-founder

sg@ltutech.com

www.LTUtech.com

- LTU has successfully deployed **image mining softwares** in **very demanding industrial environments**
 - Large data volumes, high throughput
 - Clients use the product and build value on it
 - Mission-critical tasks
 - Requirements: security, availability, failover, fast response time...and of course quality

- There are **complex issues** that need to be solved:
 - **Adapt rapidly** to changing conditions
 - Market
 - Economical, social environment
 - Technology
 - Design a **generic** and **modular** technology for multiple reuse
 - Standalone application
 - OEM
 - Integration

- Founded in 1999, LTU Technologies is a software company focused on **image mining** technologies.
 - Capitalizing on 10 years of high-profile research of the founders at **MIT Media Lab, Oxford, INRIA**
 - 20 employees, headquartered in **Paris**, office in **Washington D.C.**

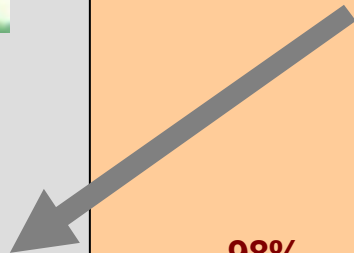
- Market verticals
 - **Law enforcement:** World-wide deployment at top-level intelligence divisions (incl. FBI).
 - *Child exploitation, Stolen Art, Counterfeiting, ID theft.*
 - **Industrial property:** Patent offices, IP-protection companies run LTU.
 - *Trademark search, Brand protection, Counterfeiting.*
 - **Asset management:** e-Commerce, publishing companies
 - *Intra/Extra/Inter-net integration to Digital Asset Management softwares.*

DNA

MD-5 or SHA-1



photo11.jpg



100%



photo11.jpg

“Duplicates”



98%

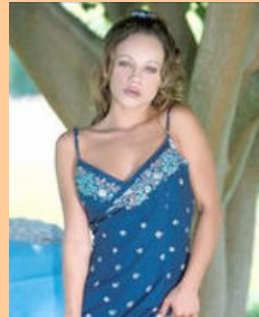


Young_girl.bmp

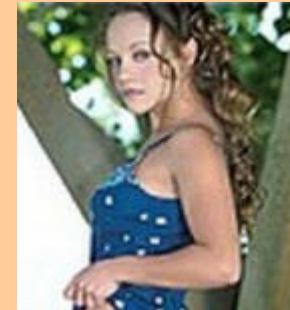
“Clones”



97%



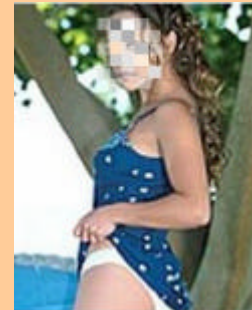
85%

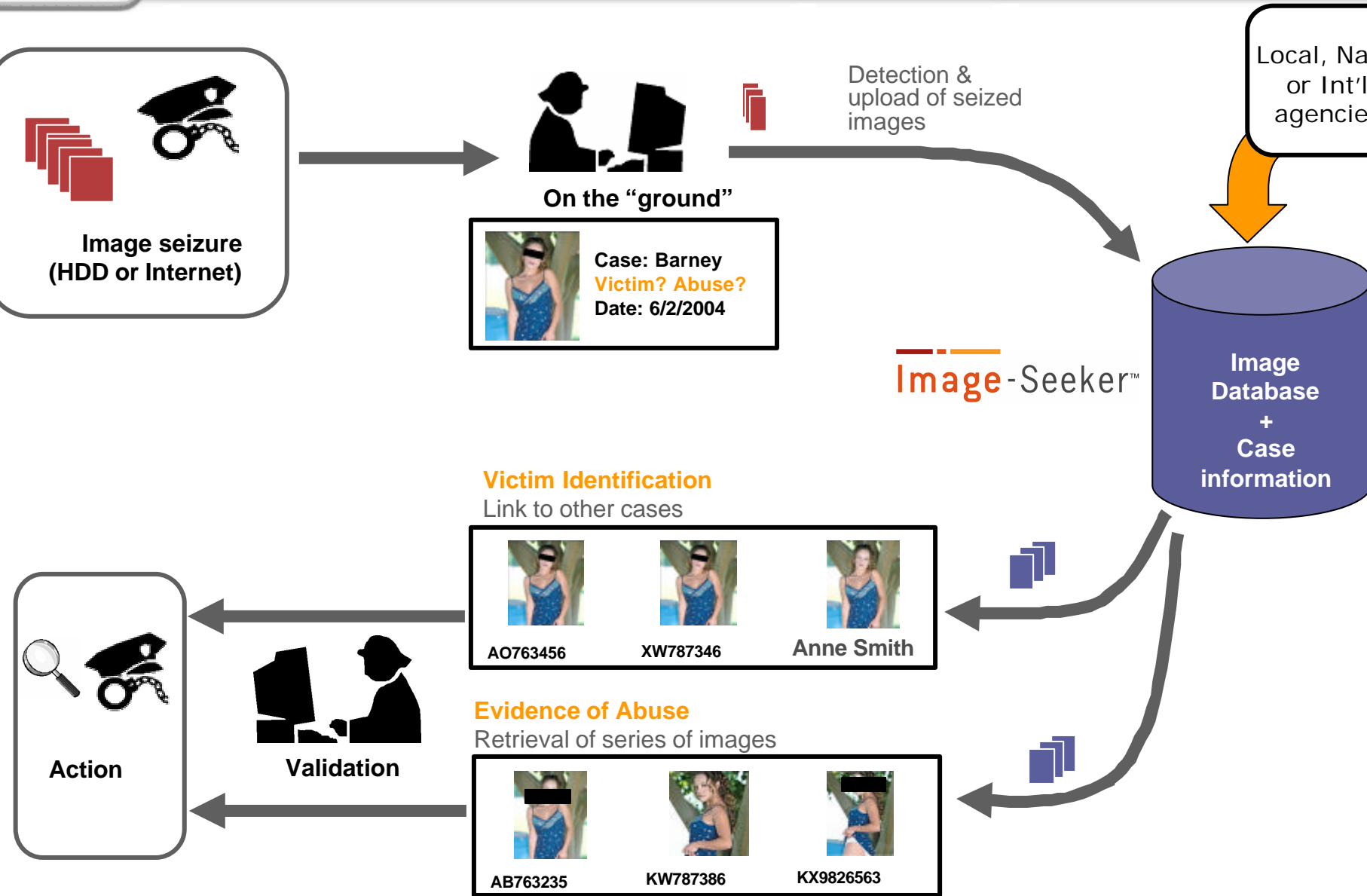


“Similar” images

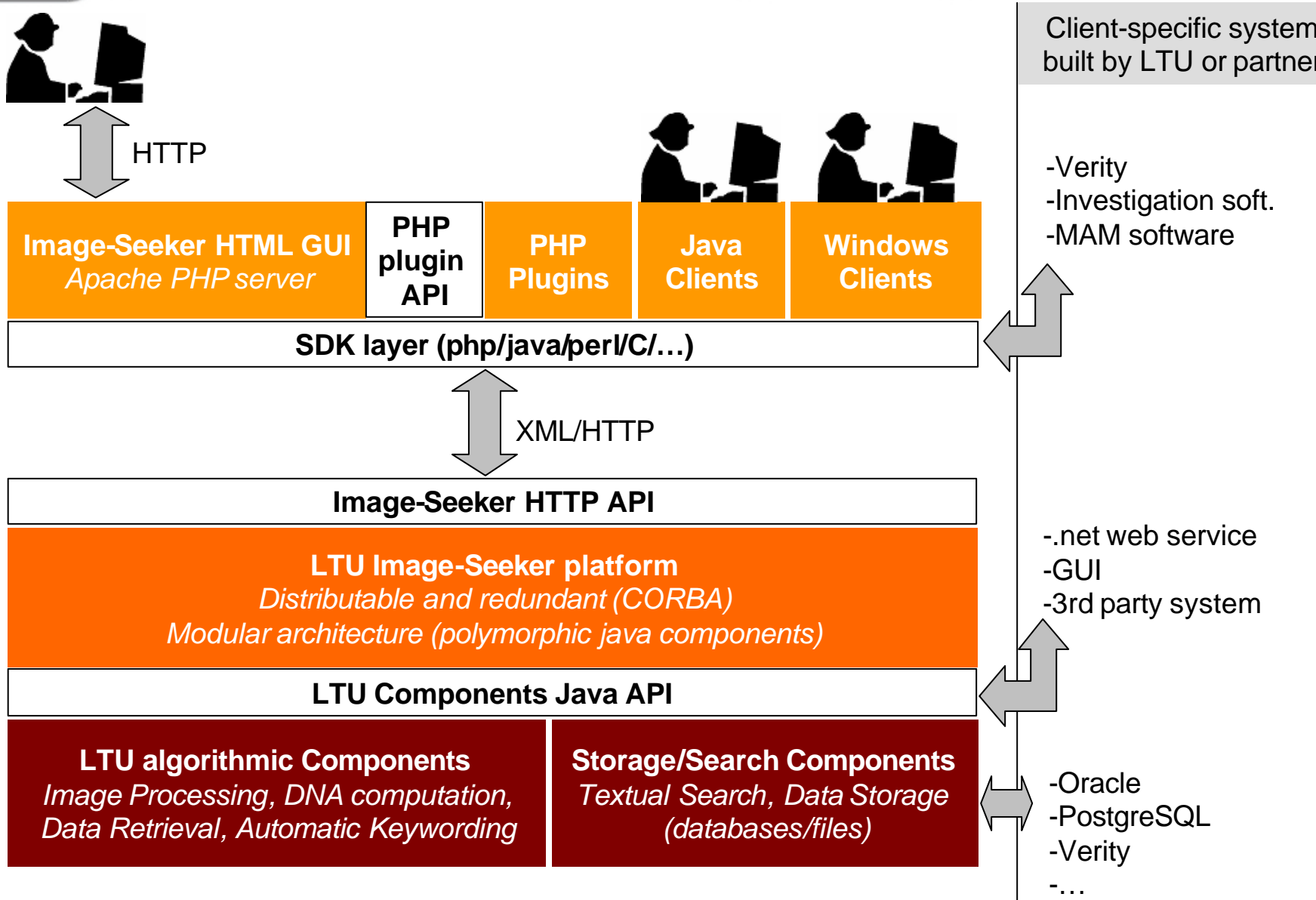


80%





A Layered Software Architecture



- General system design considerations
 - Multimedia indexing require **CPU-intensive** processes
 - Forget about DB server running MM indexing process
 - Oracle's plug-ins useless for large-scale multimedia indexing
 - N processing nodes and M storages nodes
- Ensure a true « operational » scalability
 - **Distributed, clustered architecture** (failover)
 - Bottleneck: **maintenance**
 - Re-index 1M images while maintaining QoS ?
 - Synchronization of N image repositories and data warehouses
- Adapt data models to multimedia data
 - **Increased complexity** with video, images, web pages, etc.
 - Issue: models are **task-dependent**, due to performance issues.
- Increase performance, reduce response-time
 - DBMS remain slow to access data: **need for memory caching**
 - Fast Nearest-Neighbour **search** in high dimensions and large volumes

- Real-life images are not « clean »
 - Highly **variable acquisition conditions** (uncontrolled lighting, etc.)
 - Multiple **imaging defects** (focus, over/under-exposition, etc.)
 - Collection-dependent **artefacts**

- Client requirements vary a lot
 - Heterogeneous image types (pictures, drawings, logos, etc.)
 - Highly variable definition of « **what matters in an image** »
 - Global vs. Local analysis

- Performance/Quality tradeoffs
 - **DNA extraction+classification** to be performed in **near real-time**
 - **Search** to be performed in **near real-time**
 - If asked, clients tend to favour **quality** vs. performance (insight: Moore law)

- Industry offers many **practical, technical** challenges
- Research solves many **theoretical** problems
- A gap remains, that needs to be bridged: a full client solution generally involves several technologies
 - NLP, Image Analysis, AI, etc. (e.g: trademark search)
- This means 3 main challenges !
 - **Integrating** those technologies is challenge#1
 - **Addressing real-life data**, scenarios and usages is challenge#2
 - **Optimizing** large-scale, complex systems is challenge#3
- Issues:
 - **Difficulty of obtaining** large volumes real-life **data** for academics
 - **Software re-use** too rare in academics
- Ideas:
 - **Validate research algorithms** on professional platforms... (LTU?)
 - Develop common **benchmarking** efforts on real data