# Understanding and Selecting a File Format for AV

Peter Bubestinger-Steindl

`(peter @ ArkThis.com)`

ArkThis AV-RD

# Why bother? - Let's just have:

- Best quality
- Smallest filesize
- Preserve original properties
- Fast and easy to open/use/edit
- Lasts forever
- +cherries 🍒 & ice cream 🍦on top!
- ...

# Tempting...

- Hey, it's a standard!
- Hey, everyone's using it!
- Hey, the "big ones" are using it!
- Hey, it's from a major company!
- Hey, it can do *everything*!
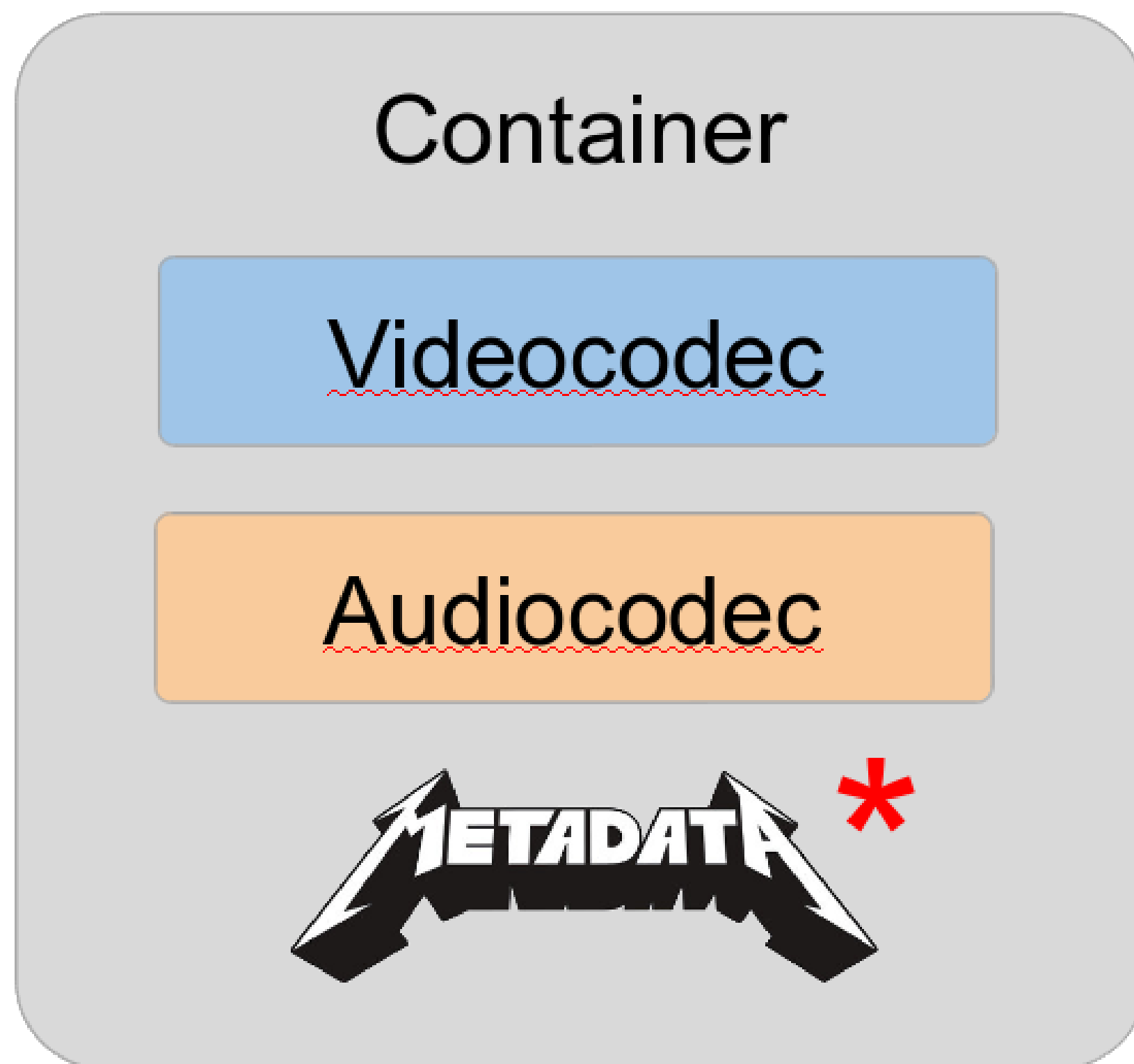- Hey, it's so easy to use!
- Hey, it's gratis!

# Digital AV formats...

What are **your** wishes, needs, expectations of a format?

# Digital AV formats...

- Which do *you* know?
- Which do *you* use?
- Which would *you* like to know more about?

# Digital Video Trinity

The uncomfortable truth. It's more than "1 format" to choose/consider for AV.

The container format is the one you see in the filename as suffix: mov, mkv, mxf, avi, flv, 3gp, wav, etc.

A/V media files: "The" format of video files usually consists of 3 different formats:

- Container
- Video codec
- Audio codec

Optional:

Embedded metadata of all kinds:

- Descriptive
- Technical
- non time-based / time-based

Additional data streams: (usually "time-based")

- subtitles
- timecode
- crazy special feature xy
- etc, etc, etc…

# What's a Container?

*"A container format (informally, sometimes called a wrapper) [...] is a file format that allows multiple data streams to be embedded into a single file, usually along with metadata for identifying and further detailing those streams."*

Source: Wikipedia: Container format (computing)

Speaker notes

Metadata: * Index (which streams are contained, etc) * Descriptive MD (title, language, etc) * as well as technical MD (fps, aspect ratio, color handling, etc) * NOTE: Some technical MD may be stored in the container and the codec/stream. This may be a blessing and a curse… Example: Aspect ratio or fps.

# What's a Container?

**Think of a regular paper folder…**

- It's a wrapper around content.
- Contains Metadata.
- Structures the content streams.

# What's a Codec?

*"A codec is a device or computer program which encodes or decodes a data stream or signal."*

Source: Wikipedia: Codec

# What's a Codec?

## Think of a human language...

- It's coded information.
- There may be dialects.
- Different people may "speak / understand" differently.

# Format Naming

Triplet notation greatly helps reducing confusion:

- **H.264** / **AAC** in **MP4**
- **FFV1** / **PCM** in **MKV** (Matroska)
- **ProRes** / **PCM** in **MOV**
- **DPX** / **WAV** (PCM) in **a folder**
- etc

# Let's look inside! :)

# VLC / MediaInfo

Website: videolan.org/vlc

Website: mediaarea.net/MediaInfo

"MediaInfo is a convenient unified display of the most relevant technical and tag data for video and audio files."

btw: A VLC-related WARNING: There's a major fraud out there: "www.vlc.de" - aka "VLC Plus Player". It contains the original VLC, but with unknown - possibly malicious - additions/modifications. Stay away from it.

# Characteristics / Properties

|  | File 1 | File 2 | File 3 |
|---|---|---|---|
| **Container** | MOV | MOV | MOV |
| **Videocodec** | UYVY | H.264 | XviD |
| **Resolution** | 720 x 576px | 1920 x 1080 | 640 x 480 |
| **FPS** | 25 | 24 | 30000/1001 |
| **-** | | | |
| **Audiocodec** | PCM | AAC | MP3 |
| **Samplerate** | 48 kHz | 48 kHz | 44.1 kHz |
| **Channels** | Stereo | Surround 5.1 | Mono |

# Significant properties

*Knowing and deciding which properties to safeguard and which are allowed to change.*

See:
LoC FADGI: DRAFT Significant Properties for Digital Video
Nestor (DE): Leitfaden DLTP AV Medien

For some it's the resolution, color information, audio quality - for others it's sufficient to see/understand it "good enough", or to be able to quickly edit-and-broadcast as the main focus.

Depends.

However, be aware that your recording may possibly be used in a different context in the future, so if possible don't aim "too low".

But please make an active decision and possibly document it (which ones and why) somewhere.

# Significant properties

Depend on media type (and use case).

| Video | Audio | Metadata |
|---|---|---|
| • resolution | • "resolution" <br> (= samplerate, bit-depth) | • language |
| • framerate | | • title |
| • aspect ratio | • channels | • author |
| • colorspace | • channel layout | • rights information |
| • subsampling | • … | • … |
| • … | | |

# "Different strokes for different folks" 😉

- **Digitization:** As-original, as-untouched as possible. Records in realtime?
  (Plus: has headroom for optional restoration/improvements.)

- **Preservation:** Stand the tests of time.
  (Highest 'original' quality)

- **Mezzanine:** For daily work. High quality.
  (Optional, if preservation format can be used for this)

- **Access** For quick and easy access.
  (Quality not necessarily best/high, but very convenient to play)

# For audio: we're lucky.

*PCM/WAV is used from digitization to preservation - and if bandwidth ain't narrow, even for access.*

Why? Because it became "small enough".

# "Different strokes for different folks" 😉

- **Digitization:** uncompressed/lossless or very-high quality lossy. (eg: V210, FFV1, MPEG-4 / PCM)

- **Preservation:** uncompressed/lossless or (very)high-quality lossy, open & documented, error-robust.
(eg: V210, FFV1, MPEG-2, MPEG-4 / PCM)

- **Mezzanine:** high-quality/high-bitrate (lossy).
(eg: MPEG-2, MPEG-4, ProRes / PCM)

- **Access** Most often lossy-compressed currently-popular format combination. (eg: H.264 / AAC in MP4)

# Format choice = A balance of ...

- Size vs Quality
- Features
- Performance
- Sustainability
- plus: time, budget, staff
- and of course: convenience

Good starting point for assessing practical usefulness.

# Risks to format longevity

- Data errors
- Obsolescence
- Vendor lock-in
- Interoperability/complexity issues

Countermeasures?

# Data errors: Error resilience?

- **Bitstream checksums:**
  Ability to *know* if the content is intact.

- **Error concealment:** Optional choice of decoder to "mask" decoding issues. (decoder specific)

- **Make backup copies!** 😇

# Open vs Closed

# Theory vs Practice



"Implementation overrules paper specs. Always."

# The Eternal Replayer





+



=

# Format Complexity

Speaker notes

Be careful with "one size fits all": Sporks are good for camping, but there's a reason why we still have separate tools for the job: spoon, knife and fork.

Considerations:

- More features = more complex / chance that only parts of specifications are supported by tool X.
- Can be non-trivial to judge what is "simple" and what is "complex"
- Find the sweet spot for your use case(s).

# Format Complexity: Less is More

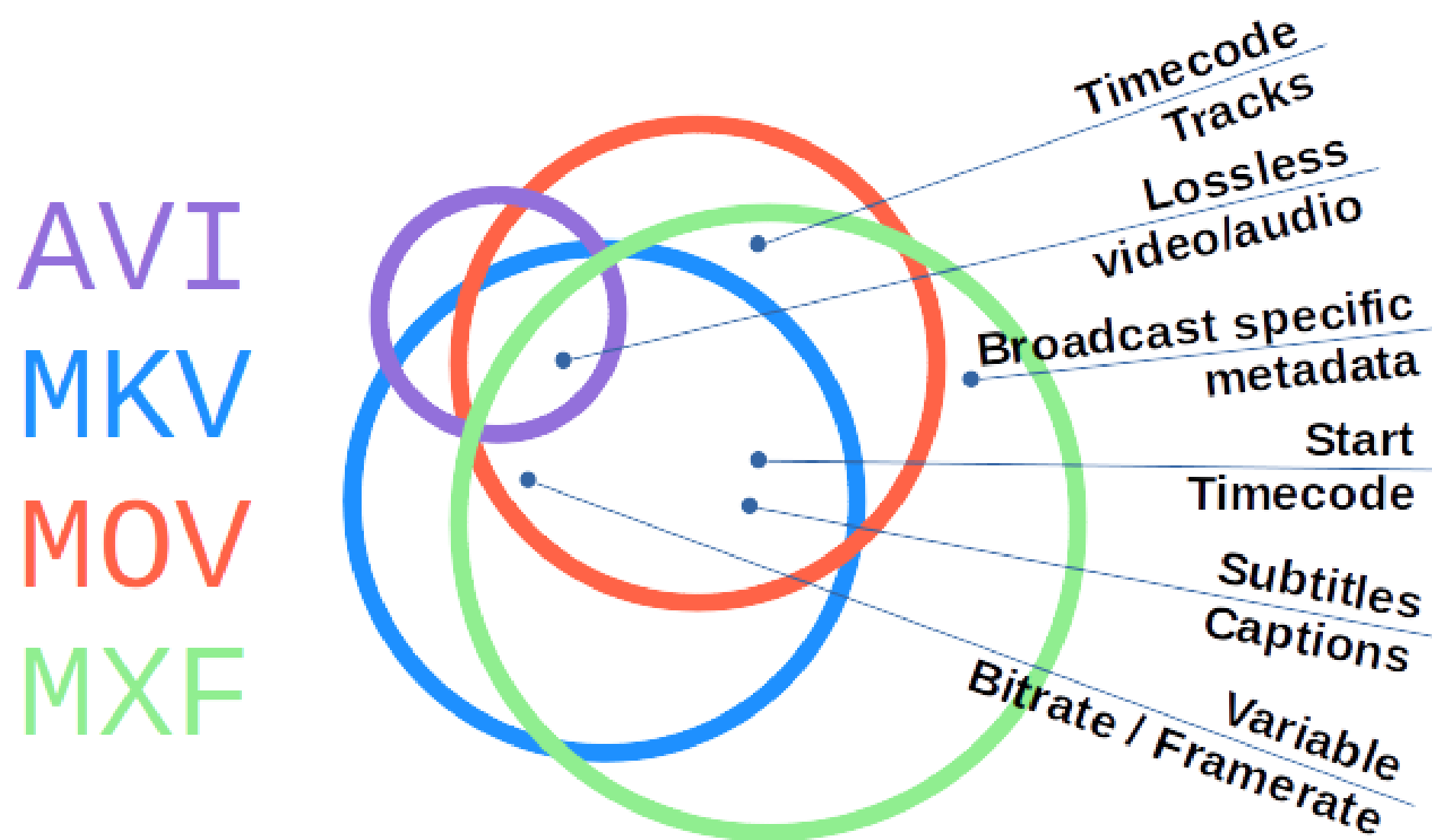Good rule = "Minimalistic Data Format":

- As simple as possible
- As complicated as necessary

*Simpler = more stable, easier to use, keep alive, reconstruct or fix.*

# Your use cases/priorities?

- Who will want/need to work with these files?
- Under which conditions?
- For how long?
- Digitization vs Production vs Preservation vs Access?
- Which properties are significant *to you*?

# Yagni Kiss Moscow?

AVI
MKV
MOV
MXF



Timecode
Tracks
Lossless
video/audio
Broadcast specific
metadata
Start
Timecode
Subtitles
Captions
Bitrate / Framerate
Variable

YAGNI / KISS / MoSCoW

# Exercise: Your Format Policy

| Must | Should | Could | Won't |
| --- | --- | --- | --- |
| _____ | _____ | _____ | _____ |
| _____ | _____ | _____ | _____ |
| _____ | _____ | _____ | _____ |
| _____ | _____ | _____ | _____ |

Choose a use-case and try to phrase your "wishes".

Now, according to the MoSCoW method, write down which features you:

- must have
- should have
- could have
- won't have (this time)

And check which format provides them, then draw a dot in the corresponding circle area.

**For example:** Only MXF may be able to provide support for broadcast-specific metadata/functions, therefore that feature will only have a dot in the MXF circle. Whereas, all (except) AVI can store aspect ratio - so that dot would go into overlap of all - except AVI. The feature of "extremely simple, well-documented and stable/unmodified for ages" would likely to be a dot in either AVI or MKV.

Use this to find out which format fits your needs, while being "as simple as possible and only as complicated as necessary".

# Examples of published Policies

- Guidelines for the Digital Film Collection
(Austrian Filmmuseum, 2018)

- Digital Preservation & Access Strategy
(Irish Film Archive, 2016)

- Digital Preservation: Policy, Standards and Procedures
(Netherlands Institute for Sound and Vision, 2016)

The above links point to preservation policies published by other AV archival institutions. They also contain information about file format choices - in the context of AV preservation.

As you can see, the length and level of detail of such policies can vary greatly. Often it is great to simply have a rather short written form that at least provides insights in "how" and "why" someone has chosen one (format) option over another.

It is also good practice to write such policies in a way that they can be interpreted and applied properly even if the actual technical options/conditions change. For example, it's okay to say "we chose format XY" - but without writing down "why", this choice may seem incomprehensible and outdated once "format XY" is superseded by another format, or something else in the tech-envirionment or the conditions has changed.

Another example: It was common to use "MP3" for some (longer) recordings in the audio archiving domain, instead of WAV/PCM. Now this may seem odd, but at that time 10 MB/Minute WAV/PCM (CD-Resolution) was heavy lifting in terms of digital computer storage.

But back then, a "large" HDD RAID was 80 GB. Consisting of several 20 GB harddisks. Now, this would fit on a USB-Stick.

*"Nothing ever doesn't change, but nothing changes much."*

All the time.

# Preservation Checklist: translated! ☆

## Sustainability:

1. Documentation openly accessible?
2. Open reference implementation?
3. How likely is it to be supported in tools/devices for which userbase?
4. Which features are implemented/tested/stable?
5. Which choice/requirements do I have to handle it beyond "shelf life"?
6. Is it legal/possible to handle it in future/different situations?
7. Can it contain proper metadata?

## Quality and functionality:

1. Preserve significant properties?
2. Sufficient image/sound quality and robustness to multi-generation copies?
3. Interoperability / ease of usage & access?
4. Direct use for editing?
5. How many different formats will I need (pile up)?
6. Handle performance / data size requirements?

# Links

- Primer on Codecs for Moving Image and Sound Archives
- Hex Editing for Archivists
- Comparing Video Codecs and Containers for Archives
- Digital Media Primer for Geeks
- A short guide to choosing a digital format for video archiving masters
- Media Digitization and Preservation Initiative (MDPI)
- Understanding audio bitrate
- Data Compression (Wikipedia)

# - Fin -

# Questions?

## Comments?

Peter Bubestinger-Steindl

`Peter @ ArkThis.com`