

BULAC

[도서관] [शिक्षक] [അക്കാദമി] [ሥልጣን]

Bibliothèque universitaire
des langues et civilisations

LA GESTION DES TRANSCRIPTIONS HTR DANS OMEKA S

Conflits d'usages au sein d'une installation complexe

Anne BUGNER · anne.bugner@bulac.fr · omeka s 4.0.4 · hébergement & administration locale sur système Ubuntu 22.04

ENS de Lyon

|

Journées Omeka - NumaHop

|

21 novembre 2024

Valoriser des projets portés par la BULAC

Projet TariMa



- Juillet 2022 - Juillet 2024
- 34 textes en arabe maghrébin (XVI^e-XIX^e siècle)
- 19 manuscrits, 9 imprimés, 6 lithographies
- 8229 pages numérisées | 7821 transcrites



Projet Chi-Know-Po Corpus



- Septembre 2022 – Janvier 2025
- 12 ouvrages de la Chine classique (III^e-XI^e siècle), éditions du XVI^e au XX^e siècle, textes glosés en colonnes
- 20+ rouleaux (8 à la BULAC)
- 33 900 pages transcrites



Usages des transcriptions dans la BiNA

- Diffusion / téléchargement des fichiers xml-alto comme médias d'un Contenu
- Visualisation en surbrillance et en texte brut du texte transcrit



- Indexation par le moteur de recherche en texte intégral

Exemple : recherche des occurrences du nom de personne “الخين لميرو”

Recherche

Ajouter le texte xml alto à la recherche en texte intégral

Activer l'option dans les Paramètres généraux

الحميد لله الى مقام العالم المشهور **الخين اميرو** معدن البضل

Saisir le texte au curseur sur l'image

Recherche

الخين اميرو

Plein texte Notice seule

Coller le texte en barre de recherche

Requête: الخين اميرو

Facettes

Appliquer Rétablir 7 résultats 1 de 1 Par page 25 Trier par Titre Liste grille

Créateur

Ibn Zaydan, Mawlā 'Abd al-Rahmān Ibn Zaydān (1873-1946 ; al-Naqib) (2)

Ibn 'Abd al-'Azīz, Muhammad (17...-1787) (2)

ابن الفتيوي, محمد (18...-1876) (1)

Auteur supposé (1)

D'après le catalogue de la fondation du roi Abdoulaye, l'auteur, Abū 'Abd Allāh Muhammad Ibn Ahmad al-Budjī al-Tarawātī al-Mansūrī al-Laknī Al-Hī, est né en 1706 et est décédé en 1775. Selon Léopold Victor Justinaud, il serait décédé en 1782 (1197 année hégirienne). Il a

7 contenus

كتاب منبع الاستماع في ذكر الجزولي والبياع وما لها من الاتباع

Dictionnaire des savants de la confrérie de cheikh al-Gazuli (m. 1465) depuis sa fondation jusqu'au XVIIIe siècle.

المنهاج السويدي

Histoire de la dynastie des Alaouites du Maroc (XVIIe-XXe siècles)

Résultats de la requête

BULAC

[도서관] [शिक्षक] [ဘာသာစကား] [ሥልጣኔ]

Bibliothèque universitaire
des langues et civilisations

Un affichage enrichi via le protocole IIIF



Affichage minimal: visionneuse Omeka S par défaut



Affichage enrichi: visionneuse IIIF Mirador

Protocole IIIF = structuration d'un ensemble d'informations relatives à 1 contenu (URL d'accès aux fichiers, résolution, description bibliographique et matérielle, conditions d'utilisation, annotation ajoutées...) dans 1 document, le **Manifest IIIF**, en vue de leur réutilisation par une visionneuse IIIF (ici **Mirador**).

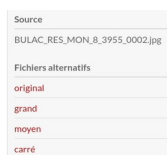
Exemple: **Contenu Omeka** et **Manifest IIIF**

BULAC

[도서관] [शिक्षक] [ဘာသာစကား] [ሥልጣኔ]

Bibliothèque universitaire
des langues et civilisations

4 modules pour 1 workflow



See also

<https://bina.bulac.fr/api/items/331388>
(application/ld+json)

3.6.18

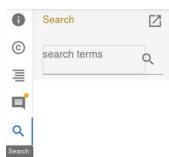
Image Server : copie du fichier source en plusieurs résolutions, assignation des fonctions (miniature, support zoom...) paramétrage avancé de l'affichage (rotation, couleur, sélection locale...), production de l'URL d'accès au fichier

IIIF manifest

<https://bina.bulac.fr/iiif/331388/manifest>

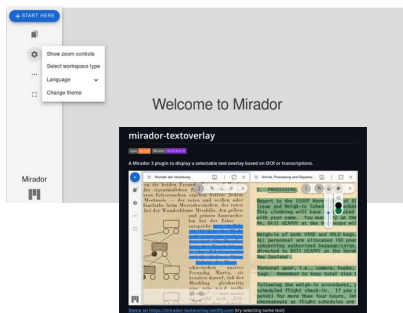
3.6.21

IIIF Server : récupération, ajout et structuration de toutes les informations nécessaires à la diffusion du Contenu et de ses fichiers média ; production du manifest IIIF



3.4.6

Iiif Search : indexation du contenu textuel du manifest IIIF dans le moteur de recherche de l'installation



3.4.9

0.3.8

Mirador Viewer : fenêtre d'affichage personnalisable du média et des informations descriptives, duplicable et ouverte au chargement de manifests extérieurs ; fonctionnalités avancées (téléchargement des ressources, annotation graphique)

Plugin **Text Overlay** de Mirador Viewer : affichage personnalisable de la transcription en surbrillance sur l'image et en volet latéral, sélection au curseur du texte

BULAC

[도서관] [शिक्षक] [ဘာသာစကား] [لغة]

Bibliothèque universitaire
des langues et civilisations

De l'xml-alto au manifest IIIF (API Presentation 2)

Section d'un fichier xml-alto : première ligne de texte de l'ouvrage

```
<Layout>  
  <Page WIDTH="982" HEIGHT="1224" PHYSICAL_IMG_NR="0" ID="page_0">  
    <PrintSpace HPOS="0" VPOS="0" WIDTH="982" HEIGHT="1224">  
      <TextBlock ID="TR_0" HPOS="93" VPOS="102" WIDTH="701" HEIGHT="916" TAGREFS="TYPE_2">  
        <Shape>  
          <Polygon POINTS="93 102 794 102 794 1018 93 1018 93 102"/>  
        </Shape>  
        <TextLine ID="TL_0" HPOS="261" VPOS="315" WIDTH="150" HEIGHT="59" BASELINE="262 363 412 362" TAGREFS="TYPE_1">  
          <Shape>  
            <Polygon POINTS="261 315 411 315 411 374 261 374 261 315"/>  
          </Shape>  
          <String ID="SE_0" CONTENT="اسنطين" HPOS="273" VPOS="315" WIDTH="131" HEIGHT="59" WC="0.9015">
```

Mapping en bref

Layout = sequence

Page = canvases 0, 1, 2, 3

PrintSpace/TextBlock = otherContent/sc:AnnotationList

TextLine = resources 0, 1, 2, 3 /cnt:ContentAsText

hpos, vpos, width, height = on

String/content = format/chars

Section du manifest IIIF

```
@context: "http://iiif.io/api/presentation/2/context.json"  
@id: "https://bina.bulac.fr/iiif/331388/annotation-page/398982/line"  
@type: "sc:AnnotationList"  
resources:  
  0:  
    @id: "https://bina.bulac.fr/iiif/331388/annotation-page/398982/line/11"  
    @type: "oa:Annotation"  
    motivation: "sc:painting"  
    resource:  
      @type: "cnt:ContentAsText"  
      format: "text/plain"  
      chars: "اسنطين"  
    on: "https://bina.bulac.fr/iiif/331388/canvas/398982#xywh=261,315,150,59"
```

Informations extraites et recomposées par le
module IIIF Server (/src/Iiif/AnnotationPage.php)

```
$xpath = $imageNumber  
? "/alto:alto/alto:Layout/alto:Page[@PHYSICAL_IMG_NR='$imageNumber']//alto:TextLine"  
: '/alto:alto/alto:Layout//alto:TextLine';
```

Bloc ciblé par le plugin Text Overlay :
resource/@type:cnt:ContentAsText + contenu

Requirements for supported IIIF manifests

- Line-level annotations with either one of:
 - a motivation that is supplementing (IIIF v3)
 - a resource that has a @type that is cnt:contentAsText (IIIF v2)

(Consignes d'utilisation sur le dépôt de code)

Objectif parallèle : pérenniser le signalement des Contenus



Utiliser des URL aussi peu susceptibles de changer que possible ==> recours à des identifiants uniques et non significants (prévenir les bugs de modules / la migration des installations / les recotations / etc)

Choix des identifiants **ARK** : gratuité, adaptabilité (niveau collection, contenu, média...) et déjà largement utilisés

Attribution d'ARK dans Omeka : module **Ark**, option **NoId** pour une génération aléatoire mais non ambiguë à la lecture humaine (exclusion des voyelles, des majuscules, du chiffre 1...)

Modification de l'URL d'accès : module **Ark Url**

NAAN	73193
NAA	example.org
Sub NAA	sub
Processeur pour le nom ark	<input type="radio"/> Id interne de la ressource <input checked="" type="radio"/> NoId
Modèle noId	bu/reeseek
Qualifiant pour le média	<input type="radio"/> Id interne du média <input checked="" type="radio"/> Position du média
Format pour la position du qualifiant	p1Id
Enregistrer l'ark avec le qualifiant pour le média	<input checked="" type="checkbox"/>
Propriété où enregistrer l'identifiant (généralement dcterms:identifier)	Dublin Core: Identifiant

Ark Url

Ark Url is a companion module for [Ark](#) or [Quark](#) that replace resources URLs by their ARK URL. So instead of `/s/site/item/1` you will get `/s/site/ark:/99999/bapZs2`.

It replaces URLs for items, item sets and media. It doesn't change anything on the admin side.

Total : 158 fichiers ajoutés à Omeka + 1 librairie

Modèle d'un identifiant ARK : ark:/ **73193** / **bmsbw4** / **p1** .thumbnail
type naan name qualifiants divers

Évolution de la syntaxe des URL BiNA :

bina.bulac.fr/s/tarima/item/331388 => bina.bulac.fr/tarima/BULAC_RES_MON 8_3955 => bina.bulac.fr/s/tarima/ark:/73193/bmsbw4

Favoriser le réemploi des manifests IIIF

- Rebond depuis les autres canaux de signalement des ressources de la BULAC : Calames, **Nakala**
Contrainte : évolution récente, mise à jour *globale* nécessaire

dcterms:requires		dcterms:URI		https://bina.bulac.fr/iiif/ark:/73193/bmsbw4/manifest	
1	id :	759			
2	fool[0]	fool[1]	fool[3]	fool[4]	fool[5]
3	Identification	Contenu des Elements et attributs à mettre à jour dans un opnlet spécifique			
4	id	ribut TITLE du 1er daoloc	l'attribut HREF du 1er daoloc type	l'eventuel daoloc au	aleur de l'attribut HREF du 1er daoloc type manifest_
5	4333	Calames-202302091524575551	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/bb0hc9	ark:/73193/bb0hc9
6	471	Calames-201644194214231	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/b9g10	ark:/73193/b9g10
7	471	Calames-201644194214232	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/b9g10	ark:/73193/b9g10
8	471	Calames-2016441942142326	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/bvdp9n	ark:/73193/bvdp9n
9	471	Calames-2016441942142327	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/b05qrv	ark:/73193/b05qrv
10	471	Calames-2016441942142328	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/b3xsc3	ark:/73193/b3xsc3
11	471	Calames-2016441942142329	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/b7pw3z	ark:/73193/b7pw3z
12	471	Calames-2016441942142330	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/bcc31k	ark:/73193/bcc31k
13	471	Calames-2016441942142331	Consultable en ligne sur la BINA	https://bina.bulac.fr/bina/ark:/73193/bcc31k	ark:/73193/bcc31k

Mention du manifest IIIF dans le dépôt des fichiers xml-alto du BULAC RES MON 8 3955 sur Nakala

Extrait du tableur ABES-BULAC pour l'ajout en masse de liens externes dans Calames

- Usages spontanés : potentiellement n'importe quel Contenu

Images

BULAC and BnF Images are available through the libraries' IIIF server. For the list of IDs (images and documents), see the [list-images.tsv](#) file. To request an image, please use the following URL template:

BULAC library (BINA)

https://bina.bulac.fr/iiif/2/{image_ID}/{region}/{size}/{rotation}/{quality}.{format}

Consigne aux utilisateurs pour retrouver les images de travail – dataset du projet TariMa sur [GitHub.com](#)

Dans tous les cas : rendre possible la déduction de l'URL du manifest depuis l'URL du Contenu

<https://bina.bulac.fr/s/tarima/ark:/73193/bmsbw4> => <https://bina.bulac.fr/iiif/ark:/73193/bmsbw4/manifest>

Intégration de l'ARK dans le manifest IIF

Module *Clean Url* : indique une syntaxe fixe pour la valeur de la propriété *dcterms:identifier*, qui sera appelée en lieu et place de l'identifiant interne par les autres scripts dans Omeka S

Modèle d'un identifiant	<input type="text" value="[a-zA-Z0-9][a-zA-Z0-9_-]*"/>
Modèle facultatif pour l'identifiant court	<input type="text"/>
Propriété pour l'identifiant	* <input type="text" value="Dublin Core : Identifiant"/>
Préfixe pour trouver l'identifiant ▶	<input type="text" value="ark:/73193/"/>
Le préfixe fait partie de l'identifiant	<input checked="" type="checkbox"/>
Les identifiants ont une barre "/" à ne pas échapper	<input checked="" type="checkbox"/>
Les identifiants sont sensibles à la casse	<input type="checkbox"/>

Total : 69 fichiers ajoutés, dont
1 au niveau racine de l'application

Paramétrage dans IIF Server

Options avancées pour les urls

Ajouter la version à l'url (à définir dans le fichier module.config.php) ▶	<input type="checkbox"/>
Utiliser les identifiants de Clean Url	<input checked="" type="checkbox"/>
Préfixe à utiliser pour l'identifiant (à indiquer dans le fichier module.config.php actuellement) ▶	<input type="text" value="ark:/73193"/>

9 scripts prévoient l'usage de Clean Url.

Modèle d'un identifiant selon ces paramètres : *prefix + name*

Interactions dans IIIF

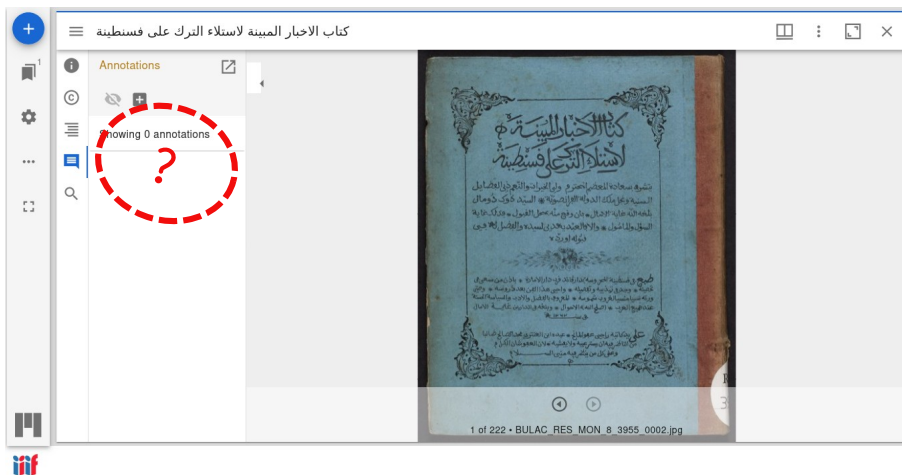
- Identifiant ARK utilisé pour l'URL du manifest (API Presentation)
- Identifiant ARK non utilisé pour l'URL image (API Image)
(... aucune modification apportée à Image Server)
- Disparitions des Annotations, et donc de Text Overlay

IIIF manifest

<https://bina.bulac.fr/iiif/bmsbw4/manifest>

See also

<https://bina.bulac.fr/api/items/331388>
(application/ld+json)



- Propriété `resource/"type": "cnt:ContentAsText"` non générée

```
status: "error"
message: 'Entité Omeka\\Entity\\Item avec critère {"id": "bmsbw4"} non trouvée'
reason: "error-controller-cannot-dispatch"
display_exceptions: true
controller: "IiifServer\\Controller\\PresentationController"
controller_class: null
```

... Où s'est bloquée l'information ?

- Base de données Omeka ? Scripts PHP ? Javascript de la visionneuse?
 - Le message d'erreur était bien sur le manifest
- Serveur virtuel ? Serveur institutionnel? Navigateur?..
- Module IIIF Server ? Clean Url ? Ark? Common?.. Ou n'importe quel autre?
- Quel impact de la version de l'API IIIF Presentation? (2 | 3)
- Incompatibilité identifiant ARK (3 parties) // Clean Url dans IIIF Server (2 parties) ?
- Méthode d'appel d'un id non-interne manquante ?...
 - Combien de fois doit-on appeler l'identifiant d'un Contenu pour IIIF? Clean Url est-il bien convoqué pour toutes?
- Les slashes ont-ils dénaturé quelque chose ?
 - Pourquoi un *name* en chaîne de caractères ne peut pas être un *name* en chaîne de caractères ?
- Que faut-il vérifier entre l'*identifier*, l'*id*, le *prefix*, le *raw identifier* et le *name* ?
 - Qu'est-ce qu'un *Controller*?
 - Les CORS?
 - Ou était-ce le *Trait* ?
- Que choisir entre la citabilité et la performance ?



```

$id = $this->params('id');
$baseAnnotationUrl = $this->view->liifUrl($item, 'iiifserver/uri', '2',
    'type' => 'annotation-list',
    'name' => $mediaId,
    'subtype' => 'annotation',
    $this->prepareMediaId());
$this->identifiersFromResources = $getIdentifiersFromResources;
$this->prefix = $prefix;
$this->rawIdentifier = $rawIdentifier;
try {
    $media = $api->read('media', ['item' => $id, 'id' => $name])->getContent();
    $array = function($xml, &$lines) use (&$xmlToArray) {
        $lines[] = [
            'name' => (string) ($xml['id'] ?? null),
            'label' => (string) ($xml['label'] ?? null),
            'children' => (string) ($xml['range_standard'] ?? null)
        ];
    };
    $iiifMediaUrl = $this->viewHelp->iiifMediaUrl($resource->id());
    $version = $this->requestedVersion();
    $params = [
        'version' => $version,
        'prefix' => empty($params['prefix']) ? $settings->get('iiifserver/uri', '2') : $params['prefix'],
        'id' => $params['id'],
    ];
    $url = $iiifMediaUrl($resource, null, $version, $params);
    if ($this->mediaIdentifier === 'media_id') {
        $element->setAttribute('id', $name);
        try {
            return $this->api()->read($url, $params);
        } catch (Exception $e) {
            return null;
        }
    }
    if ($prefixMedia) {
        $urlEncodedPrefixMedia = rawurlencode($prefixMedia);
        $constraintPrefixMedia = $prefixMedia . ' ' . $urlEncodedPrefixMedia;
        $prefixMedia = '[' . $prefixMedia . $urlEncodedPrefixMedia . ']';
    }
}
    
```

Conclusion

Tenter d'implémenter des transcriptions HTR réutilisables au sein d'un manifest IIIF à l'identifiant unique et pérenne dans Omeka S pointe **deux difficultés du modèle** :

- Une application « clés en main », prête à la personnalisation, mais dont les outils d'enrichissement se multiplient de manière non concertée ;
- Une accumulation d'extensions qui peut aboutir à une dispersion de l'information.

Comment former les gestionnaires non-développeur·euses au diagnostic des problèmes dans un environnement aussi foisonnant qu'une application Web ?

Comment améliorer la compatibilité entre extensions, sachant que les usages se traduisent souvent par des combinaisons locales de versions différentes, parfois sur des systèmes différents ?

Quelles méthodes de travail envisager au sein de la communauté de développement, au regard des autres impératifs comme la montée en performance ?