

# دليل توفير 65-80% من Claude Pro 🚀

22 نقطة عملية مع أمثلة واقعية — X Mentorship | أحمد عماد

## 1) رفع PDF خام مقابل تحويله لـ md

### المشكلة

لما بترفع PDF خام لـ Claude، كل صفحة بتأخذ ~3,000 token بسبب الصور والـ metadata. ملف 20 صفحة = 60,000 token قبل ما تكتب كلمة واحدة.

### الحل

حوّل الـ PDF لـ Markdown أو نص عادي قبل الرفع. النص الخالص بياخذ ~200 token للصفحة — **توفير فوري 93%**.

### الطريقة 1 — Google Docs (مجاني)

- افتح [drive.google.com](https://drive.google.com) وارفع الـ PDF.
- كليك يمين ← **Google Docs** ← **Open with**.
- من القائمة: **File** → **Download** → **Plain Text (.txt)** أو **Markdown (.md)**.
- ارفع الـ txt أو الـ md لـ Claude بدل الـ PDF.

### الطريقة 2 — أدوات أونلاين مجانية

- [smallpdf.com](https://smallpdf.com) → PDF to Word → انسخ النص
- [ilovepdf.com](https://ilovepdf.com) → PDF to TXT

### مثال واقعي

❌ **الغلط:** ترفع "contract\_50pages.pdf" مباشرة →  $token\ 150,000 = 3,000 \times 50$

✅ **الصحيح:** تحوّل الـ PDF لـ md وحجمه بقى 10,000 token → **توفير 93%**

⚠️ لو الـ PDF فيه جداول أو مخططات مومعة، اسكان الجزء النصي بس وابتعت الصور كـ screenshots منفصلة.

## 2) تحديد scope الملفات في Cowork

### المشكلة

Cowork يبني ملفات حسب حاجته هو — يعمل مجلدات وملفات احتياطية إنت مش طلبتها وده بياكل context.

### الحل

حدد الـ scope بوضوح في أول chat. بلّغ Claude بالضبط إيه اللي محتاجه قبل ما يبدأ يشتغل.

### Template أول رسالة في Cowork

المهمة: [اوصف المهمة]  
الملفات المطلوبة فقط:  
- [وظيفته]: [اسم الملف 1]  
- [وظيفته]: [اسم الملف 2]  
لا تعمل أي ملفات إضافية غير اللي ذكرتها.

✗ **الغلط:** "ابني لي تطبيق Claude → todo list" هيعمل 7 ملفات معظمها مش محتاجها

✓ **الصح:** حدد index.html + app.js فقط → context نظيف، tokens أقل

## 3) طول الـ Prompts

### المشكلة

ناس بتكتب prompts من 300-500 كلمة وهي مش محتاجة. كل كلمة زيادة في الـ prompt = زيادة في كل رسالة.

### الحل

29 كلمة كافية في معظم الأحيان. الـ prompt المثالي: "**عايز [المهمة] عشان [الهدف]**" وخلي Claude يسألك لو محتاج تفاصيل.

عايز [المهمة] عشان [الهدف].  
لو عندك أسئلة توضيحية، اسألني قبل ما تبدأ.

✗ **الغلط:** 450 كلمة تشرح فيها كل التفاصيل من البداية

✓ **الصح:** "عايز خطة لإدارة مشروع تقني متأخر عشان أوصله في الـ deadline. اسألني الأسئلة اللي محتاجها الأول" → توفير 94%

متى تكتب prompt طويل؟ فقط لو المهمة تقنية ومحتاجة specs محددة جدًا، أو عندك قيود صارمة من framework أو style guide.

#### 4 إرسال رسائل منفصلة لكل مهمة

##### المشكلة

لما بتبعث رسالة فيها 3 مهام منفصلة، Claude يحتاج يقرأ كل ال context لكل مهمة — ده بيضاعف استهلاك ال tokens.

##### الحل

**رسالة واحدة = مهمة واحدة.** اطلب كل حاجة في رسالة منفصلة جوا نفس ال chat.

✗ **الغلط:** "1. لخصي النص ده 2. اقتراحي عنوان 3. اعمل outline" — كل ده في رسالة واحدة

✓ **الصح:** رسالة 1: "لخصي النص ده" ← انتظر الرد ← رسالة 2: "اقتراحي عنوان" ← انتظر الرد ← رسالة 3: "اعمل outline"

في الطريقة الأولى، Claude بيحسب ال context الكامل 3 مرات. في الطريقة الثانية، كل رسالة بتأخذ context أقل — توفير 40-60%.

#### 5 تعديل الرسالة الأصلية بدل Regenerate

##### المشكلة

لما بتعمل Regenerate أو بتكتب "مش كده، اعمل تاني"، Claude بيقرأ كل ال history من الأول + رسالتك الجديدة. زيادة tokens بدون فائدة.

##### الحل

اضغط **Edit** ✏ على رسالتك الأصلية، صلح فيها، واضغط **Save & Submit**. ال history هيتبدل ومش هيتراكم.

✗ **الغلط:** أنت: "اكتبي Claude" ← email يكتب ← أنت: "مش كده" ← Claude يكتب تاني ← أنت: "لأ، باللغة العربية" — 3 رسائل + 3 ردود =  $tokens \times 3$

✓ **الصح:** اضغط Edit على رسالتك الأصلية وحددها كاملة من أول ← رسالة 1 + رد واحد = نص ال tokens — توفير 50-70%

## 6 اختيار الموديل الصح

### المشكلة

ناس بيستخدموا Opus لكل حاجة. Opus تقريبًا 5x أعلى من Sonnet. لو 70% من مواهك بسيطة، إنت بتدفع 5x زيادة.

### دليل الاختيار

الموديل	استخدمه لـ
Haiku	الأسئلة السريعة، التلخيص البسيط، الترجمة، الردود القصيرة
Sonnet	الكتابة، التحليل، الكود المتوسط، معظم المهام اليومية
Opus	التفكير المعقد، الاستراتيجية، القرارات الصعبة، الكود المتقدم
Opus + Extended Thinking	المسائل الرياضية الصعبة، التحليل العميق

كيف تغير الموديل؟ في أعلى الـ chat، اضغط على اسم الموديل الحالي واختار المناسب للمهمة — ممكن تغير في أي وقت جوا نفس الـ chat.

## 7 إلغاء الملفات الاحتياطية غير الضرورية في Cowork

### المشكلة

Cowork بيرمي ملفات "احتياطية" في كل مهمة — logs, cache files, backup versions — وكلها بتاكل tokens وبتشوش الـ context.

### الحل

حدد الملفات اللي المهمة محتاجها فعلاً في بداية الـ session.

دي session قواعد للـ

1. logs لا تعمل ملفات
2. backup files لا تعمل
3. ما لم أطلبه صراحةً README لا تعمل
4. الملفات المطلوبة فقط: [قائمة الملفات]

**النتيجة:** بدل 8 ملفات (~12,000 token)، بتأخذ 2 ملفات (~3,000 token) — توفير 75%

## 8 إدارة طول ال Chat

### المشكلة

Context في ال chat الواحد بيكبر بشكل quadratic مش linear — كل رسالة جديدة بتقرأ كل الرسائل اللي فاتت. chat بـ 40 رسالة بياخد context × 40 مقارنة بـ chat جديد.

### الحل

كل 15-20 رسالة: لخص ال brief وافتح chat جديد.

ده في 5-7 نقاط session لخصلي أهم نقط الـ جديد وأبدأ من نفس النقطة chat بحيث أقدر أفتح

بعد ما تاخذ الملخص، افتح chat جديد والصق الملخص كأول رسالة مع المهمة الجاية — توفير 50-70%.

## 9 فصل المواضيع في chats مختلفة

### المشكلة

لما بتحت 3 مواضيع مختلفة في chat واحد، Claude بيعيد قراءة كل ال context الميت وهو بيشتغل على الموضوع الجديد. tokens ضايعة.

### الحل

موضوع واحد = chat واحد.

- Chat للكود والتقنية
- Chat للكتابة والمحتوى
- Chat للتحليل والبحث
- Chat لكل project مستقل

قبل ما تبدأ أي موضوع جديد: هل الموضوع مرتبط بال chat الحالي؟ نعم → كقل | لا → افتح chat جديد — توفير 50-70%

## 10) تقليص ملف الـ about-me

### المشكلة

ملف الـ about-me بـ 22,000 كلمة بيتبعث في كل رسالة في كل session. لو عندك chat بـ 20 رسالة، الـ 22,000 كلمة بتتكرر 20 مرة.

### الحل

قلّص الـ about-me لأقل من 2,000 كلمة. وفي آخر كل session، اكتب ملاحظاتك في session-notes.md منفصل.

### الـ about-me المثالي

الاسم: [اسمك]  
الدور: [وظيفتك]  
اللغة المفضلة: [عربي/إنجليزي/مختلط]  
أسلوب الرد: [مختصر/تفصيلي/نقاط/سرد]  
سياق مهم:  
- [نقطة 1]  
- [نقطة 2]  
- [نقطة 3]

التوفير: من 22,000 × كل turn لـ 2,000 × كل turn — توفير 90%

## 11 إيقاف Search وال Connectors الافتراضية

### المشكلة

Search وال Connectors شغّالين by default. كل رسالة بتكّلف tokens إضافية لـ feature ده حتى لو مش محتاجاه.

### الحل

اوقف كل feature افتراضيًا. شغّل كل feature حسب المهمة مش حسب الـ account.

### خطوات التنفيذ في claude.ai

1. اضغط على Settings ⚙️
2. روح لـ Features أو Integrations
3. وقّف: Web Search, Google Drive Connector, Code Execution
4. شغّل أي feature بس لما تحتاجه في المهمة المحددة

شغّله لما	Feature
تحتاج معلومات بعد 2024 أو أخبار حديثة	Web Search
تشتغل على ملفات Drive مباشرة	Google Drive
تحتاج تشغّل كود وتشوف النتيجة	Code Execution

توفير: 10-20% من overhead كل رسالة

## 12 استخدام Projects بدل رفع نفس الملف مرات

### المشكلة

لو بترفع نفس الـ PDF أو الـ codebase في 5 chats مختلفة، إنت بتحرق 5 × tokens على نفس الملف.

### الحل

Projects في claude.ai بتخليك ترفع الملف مرة واحدة، وكل chat جوا الـ Project بيشوفه تلقائيًا.

### خطوات التنفيذ

1. في claude.ai، اضغط **Projects** من القائمة الجانبية
  2. اضغط **New Project** واڊي الـ Project اسم واضح
  3. ارفع الملفات الثابتة مرة واحدة (documents, codebase, style guides)
  4. ابدأ كل chat جديد من جوا الـ Project — الملفات موجودة تلقائيًا
- ✗ بدون 3 = Chat 1 + Chat 2 + Chat 3 = رفع 15,000 token على نفس الملف
- ✓ مع Project: رفع مرة واحدة فقط — توفير 80-90%

## 13 ضبط Personal Preferences مرة واحدة

### المشكلة

ناس بتكتب في كل chat: "رد بالعربي"، "استخدم bullet points"، "كن مختصر"... ده tokens ضايعة في كل مرة.

### الحل

Settings → Personal Preferences — اكتب كل تفضيلاتك مرة واحدة وهيفضلوا للأبد.

```
اللغة: رد بالعربي دائماً ما لم أطلب تحديداً غير كده
الأسلوب: مباشر ومختصر بدون مقدمات
للقوائم والنقاط المتعددة bullet points التنسيق: استخدم
الأمثلة: دائماً اضرب مثال واقعي على أي نظرية
الطول: الرد يكون بالقدر اللي المهمة تحتاجه
التعامل مع الكود: اشرح الكود بعد ما تكتبه
```

التوفير: 100% على التوضيحات المتكررة في كل chat

## 14 بناء Prompt Library

### المشكلة

بتعيد كتابة نفس ال prompts من الصفر في كل مرة. مهام زي "حل الكود ده"، "اكتبلي email احترافي"، "الخصلي الاجتماع" — كلها بتتكتب من جديد.

### الحل

اعمل ملف prompts.md وحفظ فيه ال prompts الثابتة ال templates بتاعتها. الجزء المتغير بس هو ال content اللي بتحطه.

```
## تحليل كود
:حلل الكود ده وقولي
1. المشاكل المحتملة
2. performance تحسينات ال
3. security issues أي
["الكود هنا"]

## Email احترافي
.لد [الجهة] عن [الموضوع] [ودي/رسمي] email اكتبلي
"النقاط الرئيسية: [النقاط]"
```

أو افتح Project اسمه "Prompt Library" وارفع فيه ملف ال prompts ك reference دائمة — توفير 70-80% في الوقت

## 15 استخدام Extended Thinking بحكمة NEW

### المشكلة

Extended Thinking يبأخذ tokens كثر — يفكر "بصوت عالٍ" قبل ما يجاوب. لو استخدمته على مهام بسيطة، إنت بتدفع 3-tokens زيادة بدون فائدة.

### الحل

Extended Thinking مفيد جدًا — بس في الأنواع الصح من المشاكل.

#### استخدمه لـ ✓

- مسائل رياضية أو logic معقدة
- قرارات استراتيجية بين خيارات متعددة ومتضاربة
- debugging كود صعب مش قادر تحله
- تحليل أكاديمي أو بحثي عميق

#### متستخدمهوش لـ ✗

- كتابة إبداعية أو ترجمة
- تلخيص أو أسئلة معلوماتية بسيطة
- مهام روتينية

توفير: 60-80% على المهام البسيطة بدل ما تستخدم 3-5x Opus

## 16 تقسيم المهام الكبيرة إلى Subtasks NEW

### المشكلة

لما بتعطي Claude مهمة كبيرة جدًا في رسالة واحدة، بيحتاج context ضخم جدًا ويولد output كبير. ده بيرفع الـ tokens وبتزيد احتمالية الإجابة المتوسطة.

### الحل

قسّم المهمة الكبيرة لـ subtasks صغيرة وأعطها بالتسلسل.

```
"كامل authentication المهمة الكبيرة: "ابني نظام"
```

```
Subtask 1: "database schema عمل للـ users table"
```

```
[انتظر وراجع]
```

```
Subtask 2: "POST /register endpoint عمل"
```

```
[انتظر وراجع]
```

```
Subtask 3: "POST /login endpoint عمل مع JWT"
```

توفير: 40-60% — context صغير لكل subtask بدل context ضخم + output كبير

### المشكلة

لما بتبع context كبير بدون Claude, structure بيقرأ كل حاجة بنفس الوزن. ده بيشتت انتباهه ويزود الـ tokens المطلوبة للفهم.

### الحل

استخدم XML tags لتنظيم الـ context وتحديد أولويات المعلومات.

```
<context>
  الشركة: Beetleware – SaaS solutions
  الموقع: مصر والسعودية
</context>

<task>
  Senior Python Engineer لـ job description اكتب
</task>

<requirements>
  - خبرة: +5 سنين
  - تقنيات: Python, AWS, PostgreSQL
</requirements>
```

Claude هيركز على الـ task ويستخدم الـ context كمرجع فقط — توفير 20-30%

### المشكلة

لو بنستخدم Claude API في تطبيق، كل request بيبيعت نفس الـ system prompt من جديد. لو الـ system prompt 5,000 token وعندك request 1,000 يومياً = 5,000,000 token يومياً على system prompt بس.

### الحل

استخدم Prompt Caching في الـ API. الـ cached tokens بتكلف 10% فقط من السعر العادي.

```
import anthropic
client = anthropic.Anthropic()
response = client.messages.create(
    model="claude-opus-4-5",
    max_tokens=1024,
    system=[{
        "type": "text",
        "text": "...متخصص في [مجالك] assistant أنت",
        "cache_control": {"type": "ephemeral"} # هنا السحر ←
    }],
    messages=[{"role": "user", "content": "سؤال المستخدم هنا"}]
)
```

**التوفير: 90% على الـ system prompt في كل request**

### المشكلة

معظم الناس يستخدم Projects كـ "مجلد ملفات" بس. بس الـ Project System Prompt هو أقوى أداة فيه — وناس كثير مش عارفها.

### الحل

كل Project ممكن يكون عنده **System Prompt ثابت** بيتطبق على كل chat جواه. يعني مش محتاج تكتب السياق والقواعد في كل chat من جديد.

### خطوات التنفيذ

1. افتح الـ Project من القائمة الجانبية
2. اضغط على **Project Settings** أو 
3. في خانة **Instructions** أو **System Prompt** اكتب:

#### : السياق

- المشروع: [اسم المشروع وهدفه]
- [المستخدمة stack]: التقنيات
- القيود: [أي قيود تقنية أو تجارية]

#### :قواعد الرد

- اللغة: [عربي/إنجليزي]
- الأسلوب: [مختصر/تفصيلي]

#### : افتراضات ثابتة

- [افتراض 1]
- [افتراض 2]

كل chat جوا الـ Project هيبدأ بالسياق ده تلقائيًا — مش محتاج تكتبه من جديد أبدًا.

**التوفير: 300-500 token في بداية كل chat × عدد الـ chats في الـ Project**

### المشكلة

النقطة 8 قالت "افتح chat جديد كل 15-20 رسالة" — بس إزاي تعمل الـ handoff ده هو اللي بيفرق. handoff ضعيف = بتبدأ من الصفر. handoff قوي = بتكمل بدون فقدان أي حاجة.

### الـ Handoff Template المثالي

في آخر الـ chat القديم، ابعت الرسالة دي:

```
. الجديد chat للـ Context Handoff عمل
:اكتبه في الشكل ده بالضبط

## Context Handoff
### المشروع والهدف
[جملة واحدة]
### القرارات المتخذة (مش قابلة للتغيير)
- [قرار 1]
### الوضع الحالي
[وين وصلنا بالضبط]
### الملفات/الكود الموجود
[قائمة مختصرة]
### الخطوة الجاية
[الجديد chat المهمة اللي هنبدأ بيها في الـ]
```

في أول الـ chat الجديد: الصق الـ Handoff وقول "ابدأ من الخطوة الجاية مباشرة"  
التوفير: بدل ما تنقل chat بـ 8,000 token، بتنقل Handoff بـ 400 token — توفير 95%

## 21) تعديل الـ Artifact بدل إعادة الكتابة

### المشكلة

لما Claude يعمل كود أو document كبير وإنت محتاج تعديل، كثير من الناس بيكتبوا "اعمل التعديل ده" — ف Claude بيعيد كتابة كل الـ artifact من الأول. لو الـ artifact 300 سطر كود وإنت محتاج تغيير 5 سطور، دي token مهدرة.

### الـ Pattern الصح

**بدل:** "عدّل الكود ده يعمل Claude ← "X" هيكتب 300 سطر تاني

**استخدم:**

```
"الحالي، عدّل فقط artifact في الـ"  
- [N]~ في السطر [X] اسمها function الـ -  
- غير [السلوك القديم] لـ [السلوك الجديد] -  
". لا تعيد كتابة أي حاجة ثانية"
```

### للـ Documents الـ Markdown

```
"الحالي، عدّل القسم الثالث فقط document في الـ"  
- غير الفقرة الأولى لـ [النص الجديد] -  
- جديد: [النقطة] bullet point أضف -  
". يفضل زي ما هو document باقي الـ"
```

التوفير: من إعادة كتابة 2,000-5,000 token لتعديل 200-500 token — توفير 80-90%

## 22 ال Feedback المحدد بدل "اعمل تاني"

### المشكلة

"غلط، اعمل تاني" أو "مش كده اللي عايزه" — دي أعلى 4 كلمات ممكن تكتبها. Claude هيقرأ كل ال context من الأول ويعمل output جديد كامل. ومع ذلك غالبًا النتيجة الجديدة مش أحسن لأن Claude مش عارف بالضبط إيه الغلط.

### Feedback Framework ال

- ✗ "مش كده"
- ✓ "[Z] بدل [Y] غلط - المطلوب [X] الجزء"
- ✗ "اعمل تاني"
- ✓ "احتفظ بـ [الجزء الصحيح]، عدّل بس [الجزء الغلط]"
- ✗ "مش عاجبني"
- ✓ "تعديلات محددة 3"
  1. [تعديل 1]
  2. [تعديل 2]
  3. [تعديل 3]

### مثال واقعي — Claude كتب email

Feedback صح: ✓

كويس ماعدا email الـ"  
السطر الأول - طوله شوية، بيان رسمي أكثر.  
الفقرة الثانية - اختصرها في جملتين بس.  
'استبدل' بالتوفيق' بـ 'مع خالص التحية - closing الـ  
".احتفظ بيه زي ما هو email باقي الـ"

Claude هيعمل 3 تعديلات صغيرة بدل ما يكتب email جديد كامل

التوفير: من output كامل جديد (~1,500 token) لتعديل محدد (~300 token) — توفير 80%

## جدول مقارنة التوفير في Tokens

#	النقطة	Tokens قبل	Tokens بعد	نسبة التوفير
1	رفع PDF خام مقابل تحويله لـ md	3,000 / صفحة	200 / صفحة	93%
2	تحديد scope الملفات في Cowork	8- ملفات (~12,000 token)	2- ملفات (~3,000 token)	75%
3	طول الـ Prompts	~700 token	~40 token	94%
4	إرسال رسائل منفصلة لكل مهمة	3 مهام × context كامل	context أقل لكل مهمة	40-60%
5	تعديل الرسالة الأصلية بدل Regenerate	context × عدد المحاولات	context مرة واحدة	50-70%
6	اختيار الموديل الصح	تكلفة Opus	تكلفة Sonnet	80% في التكلفة
7	إلغاء الملفات الاحتياطية غير الضرورية في Cowork	token +8,000	token 2,000	75%
8	إدارة طول الـ Chat	context quadratic	context linear	50-70%
9	فصل المواضيع في chats مختلفة	context كل المواضيع	context موضوع واحد	50-70%
10	تقليص ملف الـ about-me	22,000 × كل turn	2,000 × كل turn	90%
11	إيقاف Search والـ Connectors الافتراضية	overhead إضافي لكل رسالة	overhead 0	10-20%
12	استخدام Projects بدل رفع نفس الملف مرات	رفع × عدد الـ chats	رفع مرة واحدة	80-90%
13	ضبط Personal Preferences مرة واحدة	~50 chat / token للتوضيح	0	100% على التوضيحات
14	بناء Prompt Library	كتابة من الصفر كل مرة	نسخ template جاهز	70-80% في الوقت
15	استخدام Extended Thinking بحكمة	Opus × 3-5x على مهام بسيطة	Haiku/Sonnet	60-80%
16	تقسيم المهام الكبيرة إلى Subtasks	context ضخم + output كبير	context صغير لكل subtask	40-60%
17	استخدام XML Tags لتنظيم السياق	Claude يقرأ كل حاجة بنفس الوزن	رُكِّز على الـ task	20-30%
18	Prompt Caching للمستخدمين API	system prompt × كل request	10% تكلفة الـ cached tokens	90% على الـ system prompt

## النقاط المتقدمة (19-22)

#	النقطة	Tokens قبل	Tokens بعد	نسبة التوفير
19	Projects System Prompts	400 chat / context token	0 (محفوظ في الـ Project)	100% على context الـ setup
20	إدارة الـ Context Handoff بكفاءة	8,000 token handoff	400 token handoff	95%
21	تعديل الـ Artifact بدل إعادة الكتابة	2,000-5,000 token	200-500 token	80-90%
22	الـ Feedback المحدد بدل "اعمل ثاني"	output كامل جديد	تعديلات صغيرة	80%

